



University of Cyprus

Department of Electrical and Computer  
Engineering

Generalized Robust Estimation  
for a Class of Systems

Yiannis Socratous

Dissertation submitted as part of the requirements for the degree  
of Doctor of Philosophy at the University of Cyprus

September, 2010



Πανεπιστήμιο Κύπρου

Τμήμα Ηλεκτρολόγων Μηχανικών και  
Μηχανικών Υπολογιστών

Γενικευμένη Εύρωση Εκτίμηση  
για μια Τάξη Συστημάτων

Γιάννης Σωκράτους

Διατριβή η οποία υποβλήθηκε προς απόκτηση διδακτορικού  
τίτλου σπουδών στο Πανεπιστήμιο Κύπρου

Σεπτέμβρης, 2010

Copyright ©, by Yiannis Socratous 2010  
All Rights Reserved

Yiannis Socratous



## ABSTRACT

One of the most common signal processing tasks arising in various applications is the estimation of a signal from a noisy measurement. Most of the time imprecise *a priori* knowledge of input characteristics results in degradation of performance. This thesis presents a minimax approach to the design of robust estimation. Minimax methods are useful because they lead to constructive procedures for designing robust schemes. The main goal of the thesis is to derive robust least-square estimators for situation when the statistics or internal dynamics describing the signal and observations are not exactly known. Even though the primal focus is on robust minimax estimation, some aspect of other well-known estimation techniques, like the Maximum *A* Posteriori, the Maximum Likelihood estimation techniques, are also investigated. There are three main contributions in this thesis: 1) Modeling of uncertainty of a system using stochastic kernels, and joint distributions, and derivation of robust least-square estimators for various uncertainty sets, through a minimax approach. These uncertainty sets are defined by a Kullback-Leibler distance constraint. The results include existence of the optimal measures, and properties associated with the estimate of the true measure. Various examples, which also include MIMO communication models, are used to illustrate how the results apply to practical problems; 2) Applications of minimax theory developed to finite-dimensional autoregressive channel models in order to derive robust least-square estimators for a class of uncertain models. The methodology presented invokes a change of probability measure technique to derive recursive equations for the conditional distribution of nonlinear filtering problem. The conditional mean equation is solved explicitly to derive envelope and phase estimates for specific models such as the linear gaussian model and the non-coherent multipath model. For the non-coherent multipath model a connection with the classical least-square estimation is also presented; 3) Derivation of a generalized Maximum *A* Posteriori estimator, and a generalized Maximum Likelihood estimator are derived. The methodology used involves the introduction of an exponential function in the cost definition and the

## ABSTRACT

likelihood function of the the Maximum  $\hat{A}$  Posteriori estimation and the Maximum Likelihood estimation technique, respectively. A connection with the minimax approach is also presented and some examples are solved to illustrate the application of the results to theoretical problems.

Yiannis Socratous

## ΠΕΡΙΛΗΨΗ

Μία από τις πιο κοινές εργασίες στην επεξεργασία σήματος, η οποία προκύπτει σε διάφορες εφαρμογές, είναι η εκτίμηση ενός σήματος από μια μέτρηση που υπόκειται σε θόρυβο. Τις περισσότερες φορές ανακριβής γνώση των εκ των προτέρων χαρακτηριστικών εισαγωγής οδηγεί στην ανακριβή εκτίμηση ενός σήματος και στη μείωση της απόδοσης της σχετικής εφαρμογής. Αυτή η διατριβή παρουσιάζει μια προσέγγιση ελαχιστοποίησης -μεγιστοποίησης (minimax) στο σχεδιασμό εύρωστης εκτίμησης. Οι μέθοδοι minimax είναι πολύ χρήσιμες καθώς οδηγούν σε εποικοδομητικές διαδικασίες για το σχεδιασμό εύρωστων πλάνων. Ο κύριος στόχος της διατριβής είναι η εξαγωγή εύρωστων εκτιμητών ελαχίστων-τετραγώνων που θα μπορούν να χρησιμοποιηθούν σε περιπτώσεις όπου τα στατιστικά ή οι εσωτερικές δυναμικές που χαρακτηρίζουν ένα σήμα καθώς και οι παρατηρήσεις δεν είναι ακριβώς γνωστές. Αν και ο πρωταρχικός στόχος εστιάζεται στην εύρωστη εκτίμηση minimax, η διατριβή ερευνά και κάποιες πτυχές άλλων γνωστών τεχνικών εκτίμησης όπως είναι οι τεχνικές Μέγιστης εκ των Υστέρων (MAP) εκτίμησης και εκτίμησης Μέγιστης Πιθανοφάνειας (ML). Υπάρχουν τρεις κύριες συνεισφορές σε αυτή τη διατριβή: 1) Μοντελοποίηση της αβεβαιότητας ενός συστήματος χρησιμοποιώντας στοχαστικούς πυρήνες (stochastic kernels) και απόκοινού διανομές και εξαγωγή εύρωστων εκτιμητών ελαχίστων-τετραγώνων για διάφορα σύνολα αβεβαιότητας μέσω μιας προσέγγισης minimax. Αυτά τα σύνολα αβεβαιότητας καθορίζονται από τον περιορισμό της απόστασης Kullback-Leibler. Τα αποτελέσματα περιλαμβάνουν την ύπαρξη βέλτιστων μέτρων, καθώς και ιδιότητες που συνδέονται με την εκτίμηση του αληθινού μέτρου. Επιπρόσθετα, παρουσιάζονται διάφορα παραδείγματα, τα οποία περιλαμβάνουν και μοντέλα MIMO, όπου επιδεικνύεται η εφαρμογή των αποτελεσμάτων σε πρακτικά προβλήματα. 2) Εφαρμογή της θεωρίας minimax σε πεπερασμένα-διαστατικά (finite-dimensional) αυτοανάδρομα (autoregressive) μοντέλα καναλιών για την εξαγωγή εύρωστων εκτιμητών ελαχίστων-τετραγώνων για μια τάξη αβέβαιων μοντέλων. Η μεθοδολογία που παρουσιάζεται χρησιμοποιεί μια τεχνική αλλαγής του μέτρου πιθανότητας για την εξαγωγή αναδρομικών εξισώσεων για την υπό όρους διανομή ενός μη γραμμικού προβλήματος φιλτραρίσματος. Η υπό όρους μέση

## ΠΕΡΙΛΗΨΗ

εξίσωση λύνεται ρητά για να εξαχθούν εκτιμήσεις της περιβάλλουσας και της φάσης για συγκεκριμένα μοντέλα όπως το γραμμικό μοντέλο Gaussian και το μη-συνεκτικό (non-coherent) μοντέλο πολλαπλών διαδρομών. Επίσης, παρουσιάζεται η σχέση του μη-συνεκτικού μοντέλου πολλαπλών διαδρομών με την κλασική εκτίμηση ελαχίστων-τετραγώνων. 3) Εξαγωγή γενικευμένων εκτιμητών MAP και ML. Η μεθοδολογία που χρησιμοποιείται περιλαμβάνει την εισαγωγή μιας εκθετικής συνάρτησης στον ορισμό του κόστους και της συνάρτησης πιθανότητας των εκτιμητών MAP και ML αντίστοιχα. Τέλος, παρουσιάζεται η σχέση με την προσέγγιση minimax και επιλύονται επίσης μερικά παραδείγματα που επιδεικνύουν την εφαρμογή των αποτελεσμάτων σε θεωρητικά προβλήματα.



## ACKNOWLEDGEMENT

First and foremost, I would like to thank my academic advisor Professor Charalambos D. Charalambous, for his patience, continuous guidance and support throughout these years. I have benefited tremendously from his profound way of thinking, comprehensive knowledge, vision and enthusiasm. Also, I would like to thank all the faculty, administrative staff and fellow colleagues at the University of Cyprus.

I would like to express a special gratitude to my friends and colleagues Stelios, Michalis Markou, Ioannis, Kiriakos, Costas, Michalis Michaelides and Stojan for their help and support through past years. Furthermore, I would like to thank the student group of Professor Charalambos D. Charalambous, for their patience (especially during my presentation rehearsals) and their valuable comments.

Finally, I would like to thank my family and friends for their support. Special thanks goes to my wife Eleni, for her unconditional love, for showing patience throughout these years and for keeping me focused on my goals and what is truly important in life. Last but not least, I would like to thank my sons Marcos and Paris (he will be shortly arriving in our life) for keeping me alert and giving me an extra motive in order to complete this difficult journey.



# CONTENTS

<b>Abstract</b> . . . . .	i
<b>Περίληψη</b> . . . . .	iii
<b>Acknowledgment</b> . . . . .	v
<b>1 Introduction</b> . . . . .	1
1.1 Survey of Related Research . . . . .	4
1.1.1 Classical Estimation Methods . . . . .	4
1.1.2 Uncertainty Models . . . . .	10
1.1.3 Robust Minimax Estimation . . . . .	17
1.1.4 Relation of Mutual Information and MMSE . . . . .	23
1.2 Thesis Motivation . . . . .	25
1.3 Thesis Objective . . . . .	29
1.4 Contributions . . . . .	34
<b>2 Background Material</b> . . . . .	37
2.1 Basic Mathematical Background . . . . .	37
2.1.1 Functional Analysis . . . . .	37
2.1.2 Minimax Theory . . . . .	42
2.1.3 Measurable Space and Probability Space . . . . .	43

## CONTENTS

2.1.4	Random Variables . . . . .	46
2.1.5	Distribution Function . . . . .	49
2.1.6	Duality Relation Between KL Distance and Free Energy . . . . .	51
2.2	Change of Probability Measure . . . . .	52
2.2.1	Change of Probability Measure for Random Processes . . . . .	54
2.2.2	Change of Probability Measure for Linear Systems . . . . .	55
2.2.3	Change of Probability Measure for Nonlinear Systems . . . . .	66
2.2.4	Nonlinear Filtering Prediction and Smoothing . . . . .	71
<b>3</b>	<b>Robust Least-Square Estimation for a Class of Systems . . . . .</b>	<b>83</b>
3.1	Introduction . . . . .	83
3.2	Nonlinear Optimization . . . . .	85
3.2.1	Formulation on Abstract Spaces . . . . .	86
3.2.2	Uncertainty on the Channel Kernel and Minimax Pay-off . . . . .	88
3.2.3	Uncertainty on the $\hat{A}$ Posteriori Distribution and Minimax Pay-off . . . . .	94
3.2.4	Uncertainty on the Joint Distribution and Minimax Pay-off . . . . .	98
3.3	Examples from Estimation Theory . . . . .	101
3.3.1	Estimation of Random Variables . . . . .	102
3.3.2	Estimation of a Sequence of Random Variables . . . . .	107
3.4	Examples from MIMO Communication Systems . . . . .	109
3.4.1	Overview of MIMO Communication Systems . . . . .	109
3.4.2	Estimation from MIMO communication Systems . . . . .	110
3.5	Summary . . . . .	113

<b>4</b>	<b>Applications of Robust Estimation</b>	115
4.1	Introduction	115
4.2	The Minimax Filtering	117
4.2.1	State and Observation Models	117
4.2.2	Definition of Minimax Problem	117
4.2.3	Minimax Optimization	119
4.3	Minimax Estimation for Linear Gaussian Models	123
4.4	Non-Coherent Estimation in Multipath	126
4.4.1	Channel Model	127
4.4.2	Minimax Estimation of Phase and Envelope	128
4.4.3	Classical Estimation of Phase and Envelope	133
4.4.4	Numerical Results and Discussion	136
4.5	Summary	140
<b>5</b>	<b>Generalized MAP and ML Estimation</b>	141
5.1	Introduction	141
5.2	Generalized Maximum $\hat{A}$ Posteriori Estimation	143
5.2.1	Abstract Formulation	143
5.2.2	Derivation of Generalized Estimator	144
5.2.3	Connection with Minimax Approach	147
5.3	Generalized Maximum Likelihood Estimation	147
5.3.1	Abstract Formulation	148
5.3.2	Derivation of Generalized Estimator	148

CONTENTS

5.4	Examples . . . . .	149
5.4.1	Generalized MAP Estimator . . . . .	150
5.4.2	Generalized ML Estimator . . . . .	151
5.5	Summary . . . . .	153
<b>6</b>	<b>Conclusion</b> . . . . .	<b>155</b>
6.1	Synopsis . . . . .	155
6.2	Directions for Future Research . . . . .	157
<b>A</b>	<b>Basic Matrix Identities</b> . . . . .	<b>163</b>
<b>B</b>	<b>Proof of Remark 2.2.20</b> . . . . .	<b>165</b>
<b>C</b>	<b>Derivations of Examples of Chapter 3</b> . . . . .	<b>175</b>
<b>D</b>	<b>Proof of Theorem 4.4.4</b> . . . . .	<b>193</b>
<b>E</b>	<b>Derivations of Examples of Chapter 5</b> . . . . .	<b>197</b>
	<b>List of Acronyms</b> . . . . .	<b>203</b>
	<b>Notation and List of Symbols</b> . . . . .	<b>205</b>
	<b>Bibliography</b> . . . . .	<b>209</b>
	<b>Curriculum Vitae</b> . . . . .	<b>215</b>

## LIST OF FIGURES

1.1 Probabilistic representation of a communication channel . . . . .	11
4.1 MSE of robust minimax phase estimator for a reference SNR=20 dB .	137
4.2 MSE of robust minimax attenuation estimator for a reference SNR=20 dB	138
4.3 MSE of the phase estimator . . . . .	139
4.4 MSE of the attenuation estimator, for a reference SNR=10 dB . . . . .	139
4.5 MSE of the noiseless received signal estimator, for a reference SNR=10 dB . . . . .	140
5.1 MAP Cost Function . . . . .	145





# CHAPTER 1

## INTRODUCTION

Signal processing is an area of electrical engineering, systems engineering, and applied mathematics that deals with operations on or analysis of signals, in either discrete or continuous time to perform useful operations on those signals. Signals of interest can include sound, images, time-varying measurement values and sensor data, for example, biological data such as, electrocardiograms, control system signals, telecommunication transmission signals, such as radio signals, and many others. Signals are analog or digital representations of time-varying or spatial-varying physical quantities.

One of the most common signal processing tasks arising in applications is that of estimating (e.g., filtering, predicting, or smoothing) a signal waveform from a noisy measurement. Signal estimation is the area of study that deals with the processing of signals that contain information in order to extract this information from them. Applications of the theory of estimation is found in many areas such as communications, controls, and signal processing. This task arises, for example, in radar and sonar tracking systems, in observers for automatic control systems, in demodulators for analog communication systems, and in medical imaging systems.

In communications applications such as data transmission or radar, estimation provides the theoretical and analytical basis for the design of effective communication transmitters and receivers. Most of the time estimation applications involves taking decisions based on observations that are distorted or corrupted by noise. Moreover, the information that one wishes to extract from such observations is unknown to the observer. In such problems it is useful to formulate estimation problems in a probabilistic framework in which the unknown behavior is modeled by probability distributions.

Basic to the study of signal estimation theory is the concept of a random observation  $Y : (\Omega, \mathcal{F}) \rightarrow (\mathcal{Y}, \Sigma_{\mathcal{Y}})$ , where  $(\Omega, \mathcal{F})$  and  $(\mathcal{Y}, \Sigma_{\mathcal{Y}})$  are measurable spaces. Here  $\mathcal{Y}$  maybe a set of vectors, real numbers, or any other set. From the observation of  $Y$  one wishes to extract information for some phenomenon related to  $Y$ , and in the case of estimation problems one wishes to estimate the value of a quantity that is not observed directly. This relation between the observation and the desired information is probabilistic in the sense that the statistical behavior of  $Y$  is influenced by the value of the quantity to be estimated. This is why a model for this situation must involve a family of probability distributions on  $\mathcal{Y}$ , the members of which correspond to statistical conditions present under the various values of the quantities to be estimated. Under this model the estimation problem is to find an optimum way of processing the observation  $Y$  in order to extract the desired information. The basic features that distinguish such problem from each other are the nature of the desired information (discrete or continuous), the amount of *a priori* knowledge that is available about their desired quantities and the performance criteria by which various estimation procedures are graded. If the information we want to extract are some static parameters that do not change with time, then the problem is defined as parameter estimation. If the parameters need to be estimated are dynamic or time-varying then the problem is defined as signal estimation.

Conventional design procedures for optimum signal estimation algorithms often require an exact knowledge of the statistical behavior both of the signal of interest and of the noise corrupting the observations. For example, in the design of optimum linear estimation algorithms someone must know the spectral or autocorrelation properties of the signal and noise in order to specify the optimum procedures, and procedures designed to be optimum for a given model can be undesirably sensitive to inaccuracies or uncertainties in the model.

The two basic sources of uncertainty are noisy data and imprecise knowledge of the model of the underlying physical system. Often, *a priori* knowledge is available to describe, in a probabilistic sense, the two aforementioned sources of uncertainty, e.g., bounds, moments, mixed moments, distribution, stationarity, dependence, and spectral properties. Of course, the application of such information in the estimation procedure should yield more reliable solutions

Because the signals and noise in signal processing applications are usually modeled as random processes, performance measures usually involve probabilistic quanti-

ties (such as mean squared error or probability of error). The theory of statistics has played a fundamental role in the development of optimum signal processing techniques.

Suppose a parameter estimator for a signal with known waveform in additive noise, is designed to give optimum performance for noise possessing a specific statistical description. For example, one widespread model for noise is the Gaussian process. An important question that arises is, how sensitive is the performance of such an optimum scheme to deviations in the signal and noise characteristics from those for which the scheme is designed? This is an important question because in practice one rarely has perfect knowledge of, say, the noise characteristics; the Gaussian or any other specific model is usually a nominal assumption which may at best be approximately valid most of the time. Unfortunately, it turns out that in many cases nominally optimum signal processing schemes can suffer a drastic degradation in performance even for apparently small deviations from nominal assumptions. It is this basic observation that motivates the search for robust signal processing techniques; that is, techniques with good performance under any nominal conditions and acceptable performance for signal and noise conditions other than the nominal, which can range over the whole of allowable classes of possible characteristics. Thus, in seeking robust methods it is recognized at the outset that a single, precise characterization of signal and noise conditions is unrealistic, and so classes of possible signal and noise characterizations are constructed and considered.

Often a class of allowable characteristics, say for a noise power density function, is constructed by starting with a nominal characteristic and then including in the class all other characteristics that are “close,” in some well-defined sense, to this nominal one. Then a signal processing scheme that is robust may have performance at the nominal which is not quite as good as the scheme that is optimum for the nominal case, but its overall performance with respect to the defined class of characteristics will be good or acceptable. This loose definition of robustness is perfectly reasonable, but it does not provide a systematic approach to obtaining robust schemes. In order for this to be achieved first a measure of “overall” performance of a scheme with respect to a class of allowable conditions at the input is specified. One such measure that has been widely used and which leads to interesting and useful results in many situations is the worst case performance of a scheme over a class of input conditions. Clearly, if its worst case performance is good, then it may be concluded that a given scheme is robust. On the other hand, such a robust scheme can be found by looking

for the scheme that optimizes worst case performance. This approach leads to what are known as minimax robust schemes. A scheme that minimizes the maximum possible value of a loss function is called minimax; if performance is measured by a gain function then a maximin scheme would be sought. The term minimax is used as a general description for such schemes in all cases. Implicit in the association of minimax schemes with robust schemes is the expectation that the worst case performance of a minimax scheme will be acceptably good, being the best that can be achieved. Another expectation one has in defining robust schemes in this way is that at any nominal operating point, the performance of the minimax scheme will not be very far below that of the nominally optimum scheme, which on the other hand will have much poorer performance away from the nominal point.

The rest of the chapter is organized as follows. In Section 1.1, the review of related literature is given. In Section 1.2, the thesis goals are introduced and motivated. In Section 1.3, the statement of the problems are presented. In Section 1.4, the main contributions of the thesis are outlined.

## **1.1 Survey of Related Research**

This section summarizes the models and results on classical estimation methods and robust minimax estimation techniques which are relevant to this thesis.

### **1.1.1 Classical Estimation Methods**

In many situations arising in practice one is interested in making a choice among a continuum of possible states of nature. In particular, given a family of distributions on the observation space indexed by a parameter or set of parameters, the true value of the parameter has to be determined as precisely as possible from the observations. Such problems are known as parameter estimation problems. In this thesis two basic approaches to parameter estimation are presented, first the Bayesian, in which the parameter is assumed to be a random quantity related statistically to the observation, and a second, the Maximum Likelihood estimation, in which the parameter is assumed to be unknown but without being characterized with a probabilistic

structure.

### *Bayesian Estimation*

In estimation theory and decision theory, a Bayes estimator or a Bayes rule is an estimator or decision rule that minimizes the posterior expected value of a cost function. Suppose an unknown Random Variable (RV)  $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \Sigma_{\mathcal{X}})$  is known to have a prior distribution  $P_X$ . Let  $\Phi(Y)$  be an estimator of  $X$  (based on the measurements of a RV,  $Y$ ), and let  $C(X, \Phi(Y))$  be a cost function. The Bayes risk of  $\Phi(Y)$  is defined as  $E[C(X, \Phi(Y))]$ , where the expectation is taken over the probability distribution of  $X$ : this defines the risk function as a function of  $\Phi(Y)$ . An estimator  $\Phi(Y)$  is said to be a Bayes estimator if it minimizes the Bayes risk among all estimators.

One of the most popular Bayesian estimation techniques is the Least-Square estimation, also known as Minimum-Mean-Squared-Error estimation. In classical Least-Square estimation of RV's one is interested in finding the best estimate,  $\Phi^*(Y)$ , of a RV  $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \Sigma_{\mathcal{X}})$  from the measurements of a RV  $Y$ , by minimizing the expected value of the least-square cost function  $C(X, \Phi) \triangleq \|X - \Phi(Y)\|_{\mathcal{X}}^2$  over all functions  $\Phi : \mathcal{Y} \rightarrow \mathcal{X}, Y \mapsto \Phi(Y)$ , which is a function of  $Y$  ( $\|x\|_{\mathcal{X}}^2$  denotes the norm on  $\mathcal{X}$ ). By the orthogonal projection theorem, the solution is given by the conditional expectation:

$$\Phi^*(Y) = E[X|Y] = \int_{\mathcal{X}} x dP_{X|Y}(x|y) \quad (1.1)$$

where  $P_{X|Y}$  is the conditional distribution of  $X$  given  $Y$  [1], [2], [3].

The solution to this estimation problem is stated in terms of the *a posteriori* distribution function  $P_{X|Y}$ . This distribution contains all information available to estimate  $X$  or any nonlinear function of it. The objective is thus to find, recursively in time, the evaluation equation of the *a posteriori* distribution and to solve it (for the case of random processes). However, it is only in a few special cases that this distribution can be solved explicitly by parameterizing it using a finite number of statistics. A well known example is the case with linear dynamics and observations in additive Gaussian noise. In this case all densities involved are Gaussian, and hence the conditional distribution can be parameterized using the corresponding mean and covariance. The equations of the finite statistics are given by the Kalman Filter [1], [4], which are going to be explained in more details in the following paragraphs.

For the case of nonlinear systems, there are difficulties in obtaining the solution of the  $\hat{a}$  posteriori distribution recursion in closed form. Often, approximations have to be made to find sub-optimal nonlinear estimators. The standard method is to use the Taylor series expansion and apply linear filtering theory, giving rise to the so-called Extended Kalman Filter [5], [6], [7]. Other more sophisticated sub-optimal estimation techniques are available, e.g. reiteration, higher order filters, and statistical methods [5].

Another cost function that is sometimes applied is the absolute error, given by  $C(X, \Phi(Y)) = |X - \Phi(Y)|$ ,  $(X, \Phi(Y)) \in \mathbb{R}^2$ . The Bayes risk here is  $E[|\Phi(Y) - X|]$ , a quantity known as the mean-absolute-error, so the corresponding Bayes estimate is known as the Minimum-Mean-Absolute-Error (MMAE) estimate. The Bayes estimate in this case, denoted by  $\Phi_{ABS}(y)$ , is any point such that

$$\begin{aligned} P(X < t|Y = y) &\leq P(X > t|Y = y), \quad t < \Phi_{ABS}(y) \\ P(X < t|Y = y) &\geq P(X > t|Y = y), \quad t > \Phi_{ABS}(y) \end{aligned} \quad (1.2)$$

Note that a point  $\Phi_{ABS}(y)$  satisfying (1.2) is a median of the conditional distribution of  $X$  given  $Y = y$ . Thus the MMAE estimate is a conditional median estimate. This estimate coincides with the Least-Square (LS) estimate only when the distribution of  $X$  given  $Y = y$  has the same value for the mean and median.

Moreover, another estimation method that, although not properly a Bayes estimate, fits within the Bayesian framework is Maximum  $\hat{A}$  Posteriori (MAP) probability estimation. In Bayesian statistics, a MAP estimate is a mode of the posterior distribution. The MAP can be used to obtain a point estimate of an unobserved quantity on the basis of empirical data. It is closely related to Maximum Likelihood estimation, which is presented in the following paragraphs, but employs an augmented optimization objective which incorporates a prior distribution over the quantity one wants to estimate. MAP estimation can therefore be seen as a regularization of Maximum Likelihood (ML) estimation.

It is assumed that an unobserved RV  $X$  has to be estimated on the basis of observations  $Y$ . Applying a uniform cost criterion, leads to a procedure for estimating  $X$  as that value which maximizes the  $\hat{a}$  posteriori (discrete or continuous) density  $\eta(y, dx)$  (also referred in this dissertation as stochastic kernel). Therefore,

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}} \eta(y, dx) \quad (1.3)$$

In modeling a given statistical situation usually someone starts with the family of conditional distributions (or stochastic kernels) of  $Y$  given  $X = x$ , and for the Bayesian formulation with a prior distribution for  $X$  also. The conditional distribution of  $X$  given  $Y$  (or the stochastic kernel  $\eta(y, dx)$ ) is obtained from the prior and the conditional of  $Y$  given  $X$  (or the stochastic kernel  $\mu(x, dy)$ ) by applying Bayes' formula.

$$\eta(y, dx) = \frac{\mu(x, dy)dP_X(x)}{\int_{\mathcal{X}} \mu(x, dy)dP_X(x)} = \frac{\mu(x, dy)dP_X(x)}{dP_Y(y)} \quad (1.4)$$

The MAP estimator can be obtained using (1.4) but without the computation of the denominator of the posterior distribution ( $dP_Y(y)$ ) since this term does not depend on  $x$  and therefore plays no role in the optimization. That is,

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}} \mu(x, dy)dP_X(x) \quad (1.5)$$

The above MAP estimate of  $x$  coincides with the ML estimate when the prior  $dP_X(x)$  is uniform (that is, a constant function). The MAP estimate is a limit of Bayes estimators under a sequence of 0 – 1 cost functions, but not a Bayes estimator per se.

#### *Maximum Likelihood Estimation*

For many observation models arising in practice, it is not possible to apply the above results either because of intractability of the required analysis or because of the lack of useful complete sufficient statistic. For such models, an alternative method for seeking good estimators is needed. One very commonly used method of designing estimators is the Maximum Likelihood (ML) method. Maximum Likelihood estimation is a popular statistical method used for fitting a statistical model to data, and providing estimates for the model's parameters.

The ML estimator it is one of the most fundamental estimators in statistics. For a fixed set of data and underlying probability model, ML picks the values of the model parameters that make the data "more likely" than any other value of the parameters would make them. ML estimation gives a unique and easy way to determine the solution in the case of the normal distribution and many other problems, although in very complex problems this may not be the case. If a uniform prior distribution is assumed over the parameters, the ML estimate coincides with the most probable values thereof.

Consider a family of probability distributions parameterized by an unknown parameter  $x$  (which could be vector-valued), associated with either a known probability density function (continuous distribution) or a known probability mass function (discrete distribution), denoted as  $\mu(x, \cdot)$ . A sample  $\{y_1, y_2, \dots, y_n\}$  of  $n$  values is drawn from this distribution, and then  $\mu(x, \cdot)$  is used for computing the (multivariate) probability density associated with the observed data,  $\mu(x, dy)$ . As a function of  $x$  with  $y_1, \dots, y_n$  fixed, this is the likelihood function

$$\mathcal{L}(x) = \mu(x, dy).$$

The method of maximum likelihood estimates  $x$  by finding the value of  $x$  that maximizes  $\mathcal{L}(x)$ . This is the Maximum Likelihood estimator of  $x$ :

$$\hat{x}_{ML}(y) = \arg \max_x \mathcal{L}(x). \quad (1.6)$$

### *Kalman Filtering*

All the estimation methods discussed so far are used for designing estimators for static parameters, that is, for parameters that are not changing with time, although these methods are applied to sequences. In many applications the related problems of estimating dynamic or time-varying parameters is of interest. In the traditional terminology, a dynamic parameter is usually called a signal, so the above problem is known as signal estimation. The dynamic nature of the parameter in signal estimation problems adds a new dimension to the statistical modeling of the problems where the dynamic properties of the signal must be modeled statistically in order to obtain meaningful signal estimation procedures. Unlike the static case, an estimator of a signal is not expected to be perfect as the number of observations becomes infinite because of the time variation in the signal.

Kalman filtering provides a very useful algorithm for estimating signals that are generated by finite-dimensional linear dynamically models. The Kalman filter is an algorithm in control theory introduced by Kalman (1960) and refined by Kalman and Bucy (1961). It is a minimum-mean-square-error estimator and it is an algorithm which makes optimal use of imprecise data on a linear (or nearly linear) system with Gaussian errors, and continuously update the best estimate of the system's current state.

Kalman filter theory is based on a state-space approach in which a state equation models the dynamics of the signal generation process and an observation equation



models the noisy and distorted observation signal. For a signal  $x_k$  and noisy observation  $y_k$ , equations describing the state process model and the observation model are defined as

$$\begin{aligned}x_k &= Ax_{k-1} + w_{k-1}, \quad x_0 \sim N(Ex_0; cov(x_0)) \\y_k &= Hx_k + n_k,\end{aligned}\tag{1.7}$$

where,  $x_k$  is the  $p$ -dimensional signal vector, or the state parameter, at time  $k$ ,  $A$  is a  $p \times p$  dimensional state transition matrix that relates the states of the process at times  $k - 1$  and  $k$ ,  $w_k$  (process noise) is the  $p$ -dimensional uncorrelated input excitation vector of the state equation.  $w_k$  is assumed to be a normal (Gaussian) process  $p(w_k) \sim N(0, Q)$ ,  $Q$  being the  $p \times p$  covariance matrix of  $w_k$  or process noise covariance.  $y_k$  is the  $M$  dimensional noisy observation vector,  $H$  is a  $M \times p$  dimensional matrix  $n_k$  is the  $M$ -dimensional noise vector, also known as measurement noise,  $n_k$  is assumed to have a normal distribution  $p(n_k) \sim N(0, R)$  and  $R$  is the  $M \times M$  covariance matrix of  $n_k$  (measurement noise covariance). Often,  $(n_k, w_k, x_0)$  are assumed independent.

Here,  $\hat{x}_{k|k-1}$  is defined as the  $\hat{a}$  priori estimate (prediction) at step  $k$  from the previous trajectory of  $x$  given measurements  $y_0, \dots, y_{k-1}$ , and  $\hat{x}_{k|k}$  as the  $\hat{a}$  posteriori state estimate at step  $k$  given measurements  $y_0, \dots, y_k$ . Note that  $\hat{x}_{k|k-1}$  is a prediction of the value of  $x_k$  which is based on the previous values and not on the current observation at time  $k$ .  $\hat{x}_{k|k}$  on the other hand, uses the information in the current observation and previous  $y_0, \dots, y_{k-1}$ .

The Kalman filter is often derived by beginning with the goal of finding an equation that computes an  $\hat{a}$  posteriori state estimate as a linear combination of an  $\hat{a}$  priori estimate (prediction) and a weighted difference between an actual measurement and a measurement prediction (innovation). Hence, each estimate consists of a fraction which is predictable from the previous values and does not contain new information and a fraction that contains the new information extracted from the observation.

$$\hat{x}_{k|k} = \hat{x}_{k|k-1} + K_k \left( y_k - H\hat{x}_{k|k-1} \right).\tag{1.8}$$

The difference  $y_k - H\hat{x}_{k|k-1}$  in (1.8) is called the measurement innovation. The innovation reflects the discrepancy between the predicted value and the actual measurement. The  $P \times M$  matrix,  $K_k$ , in (1.8) is chosen to be the gain or blending factor that minimizes the  $\hat{a}$  posteriori error covariance. One form of the resulting  $K_k$  is

given by

$$K_k = \Sigma_{k|k-1} H^T \left( H \Sigma_{k|k-1} H^T + R \right)^{-1} \quad (1.9)$$

where  $\Sigma_{k|k-1}$  is the covariance of the prediction error,  $x_k - \hat{x}_{k|k-1}$ , conditioned on  $y_0^{k-1}$ .

The Kalman filter estimates a process by using a form of feedback control: the filter estimates the process state at some time and then obtains feedback in the form of (noisy) measurements. As such, the equations for the Kalman filter fall into two groups: time update equations (prediction) and measurement update equations (correction). The time update equations are responsible for projecting forward (in time) the current state and error covariance estimates to obtain the  $\hat{a}$  priori estimates for the next time step. The measurement update equations are responsible for the feedback i.e. for incorporating a new measurement into the  $\hat{a}$  priori estimate to obtain an improved  $\hat{a}$  posteriori estimate.

The time update equations can also be thought of as predictor equations, while the measurement update equations can be thought of as corrector equations. Indeed the final estimation algorithm resembles that of a predictor-corrector algorithm for solving numerical problems as shown below.

Time update (predict)

$$\hat{x}_{k|k-1} = A \hat{x}_{k-1|k-1}, \quad \Sigma_{k|k-1} = A \Sigma_{k-1|k-1} A^T + Q, \quad \Sigma_{0|-1} = \text{cov}(x_0), \quad (1.10)$$

where  $\Sigma_{k|k} = \Sigma_{k|k-1} - \Sigma_{k|k-1} H^T \left( H \Sigma_{k|k-1} H^T + R \right)^{-1} H \Sigma_{k|k-1}$ .

Measurement update (correct)

$$\begin{aligned} K_k &= \Sigma_{k|k-1} H^T \left( H \Sigma_{k|k-1} H^T + R \right)^{-1} \\ \hat{x}_{k|k} &= \hat{x}_{k|k-1} + K_k \left( y_k - H \hat{x}_{k|k-1} \right), \quad x_{0|-1} = E x_0. \end{aligned} \quad (1.11)$$

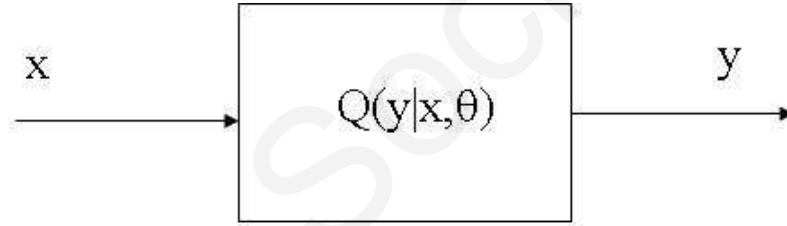
### 1.1.2 Uncertainty Models

In practical applications one of the important issues in estimation problems is the choice of an appropriate system or channel model. When the channel is uncertain, this is not a trivial problem. The uncertainty models can be divided into two

categories, the parametric uncertainty models and the nonparametric uncertainty model. In the following paragraphs some typical examples from the two categories are presented.

### *Parametric Uncertainty Models*

In the estimation theory literature a basic model, which is often used to describe a parametric uncertainty in the channel is the one depicted in Fig. 1.1 [8]. A variable  $\theta$ , which belongs to a certain set  $\Theta$ , parameterizes a conditional distribution  $Q(y|x, \theta)$ . Hence, instead of dealing with a fixed known system, one considers the estimation for the class of systems  $\{Q(y|x, \theta) : \theta \in \Theta\}$ . Generally speaking, there are two classes of uncertainty models; the class of compound systems, and the class of arbitrarily varying systems. Compound and Arbitrarily Varying Channels are further classified into discrete memoryless and finite-state systems [8].



**Figure 1.1:** Probabilistic representation of a communication channel

A family of discrete memoryless systems

$$\{Q(y|x, \theta) : x \in \mathcal{X}^n, y \in \mathcal{Y}^n, \theta \in \Theta\}_{n=1}^{\infty} \quad (1.12)$$

where

$$Q(y|x, \theta) = \prod_{t=1}^n Q(y_t|x_t, \theta_t), \quad (1.13)$$

and  $\{Q(y_t|x_t, \theta_t) : x_t \in \mathcal{X}, y_t \in \mathcal{Y}, \theta_t \in \Theta\}$  is a suitable subset of the set of all stochastic matrices  $\mathcal{X} \times \Theta \mapsto \mathcal{Y}$ , is called a discrete memoryless compound system. Thus, compound systems assume that the true system  $Q_{true}(y|x, \theta)$  is unknown, though, it is assumed that  $Q_{true}(y|x, \theta)$  belongs to the family of systems (1.12) and remains unchanged during the course of a transmission.

Arbitrarily Varying Channels (AVCs) are a generalization of (1.12) to include time variations in which the system changes during subsequent transmitted letters  $x_k$ .

Assume that  $\Sigma$  is a finite set of system states and  $\theta = \Sigma^\infty$ . Then an AVC is determined by

$$Q(y|x, s) = \prod_{t=1}^n Q(y_t|x_t, s_t), \quad (1.14)$$

where  $s = (s_1, \dots, s_n)$ , and  $\{Q(y_t|x_t, s_t) : x_t \in \mathcal{X}, y_t \in \mathcal{Y}, s_t \in \Sigma\}$  is a suitable subset of the set of all stochastic matrices  $\mathcal{X} \times \Sigma \mapsto \mathcal{Y}$ . Hence, at each moment  $t$ , the transmission matrix  $Q(y_t|x_t, s_t)$  is unknown, and it is determined by the system state  $s_t \in \Sigma$ .

Continuous alphabet uncertainty systems received much less attention in the literature than the discrete counterparts. Most of the results are related to Gaussian uncertainty systems, which are briefly described below.

A Gaussian compound channel is defined by

$$Y = HX + W, \quad (1.15)$$

The uncertainty is introduced by assuming that a linear transformation  $H$  is unknown, although it is known to belong to some pre-specified class of linear transformations, e.g.  $\|H\| \leq \gamma$ , or  $H = \tilde{H} + \delta H$ ,  $\|\delta H\| \leq p$  and  $\tilde{H}$  is known.  $X$  and  $Y$  are in general continuous signals, and  $W$  is an additive Gaussian noise. These models can also be used when  $H$  is known and there is uncertainty in  $X$ , which can be described with a constraint on the weighted norm  $\|X\| \leq L$  or a constraint on the covariance of  $X$ ,  $C_X \leq S$ .

A generalized parametric uncertainty model is being used in [9] which accounts for uncertainties in the data  $\{A, b\}$ . The uncertainties in these data, expressed as  $\{\delta A, \delta b\}$ , are assumed to lie within certain balls of radii  $\{\eta, \eta_b\}$ , i.e., they are known to be bounded and satisfy  $\|\delta A\| \leq \eta$ ,  $\|\delta b\| \leq \eta_b$ . A special case is also given in [9] where the perturbations  $\{\delta A, \delta b\}$  are assumed to satisfy a model of the form

$$[\delta A \ \delta b] = HS[E_a \ E_b], \quad (1.16)$$

where  $S$  is an arbitrary contraction,  $\|S\| \leq 1$ , and  $\{H, E_a, E_b\}$  are known quantities of appropriate dimensions. As a brief motivation, one application of state-space estimation is succinctly described in [9], with full details provided in [10]. For this case the following uncertainty model is considered

$$\begin{aligned} X_{i+1} &= (F_i + \delta F_i)X_i + (G_i + \delta G_i)U_i \quad i \geq 0, \\ Y_i &= H_i + V_i \end{aligned} \quad (1.17)$$

where the perturbations in  $\{F_i, G_i\}$  are modeled as

$$[\delta F_i \ \delta G_i] = M_i \Delta_i [E_{f,i} \ E_{g,i}] \quad (1.18)$$

for some known matrices  $\{M_i, E_{f,i}, E_{g,i}\}$  and for an arbitrary contraction  $\Delta_i$ ,  $\|\Delta_i\| \leq 1$ .

Of special interest are Multiple-Input Multiple-Output (MIMO) channel models, due to the applications of multiple-antenna systems for wireless communication. Initially, multiple-antenna systems promised considerable gain for fading channels as compared to single-antenna systems. However, at a later stage, it has been shown that this gain depends on the level of knowledge that the transmitter and the receiver have about the channel. The received signal of a flat fading channel is given by

$$y = HX + W \quad (1.19)$$

where  $x$  is a transmitted vector in  $\mathbb{C}^n$ ,  $w$  and  $y$  are random variables in  $\mathbb{C}^d$ , and  $H$  is a channel matrix in  $\mathbb{C}^{d \times n}$ . The additive noise  $w$  and channel matrix  $H$  are ergodic and stationary, and their entries are i.i.d., zero mean, circularly symmetric complex Gaussian random variables. If the channel matrix  $H$  is not perfectly known to the transmitter and /or receiver, the uncertainty is modeled as additive.

$$H = \hat{H} + E. \quad (1.20)$$

Here,  $\hat{H}$  represents the estimation of  $H$ , while  $E$  is an estimation error.

An example of a flat-fading MIMO uncertainty model is given in [11]. Given a channel matrix  $H$ , the noisy observation can be expressed as (1.19) where for this problem  $x$  is the  $\eta_T$ -dimensional symbol vector transmitted during a signaling period and denotes a zero-mean complex circular Gaussian noise vector with covariance matrix  $R$ . It is assumed that  $R$  is positive definite so that the noise affects all observation components, and  $x$  has zero-mean and normalized covariance matrix  $I_{\eta_T}$  and is independent of  $w$ . This paper considers a robust Mean-Square-Error (MSE) equalizer design for MIMO communication systems with imperfect channel and noise information at the receiver. When the estimated channel  $\hat{H}$  and the noise covariance matrix  $\hat{R}$  are different from the actual channel  $H$  and noise covariance matrix  $R$ , respectively, the MSE objective function is expressed as a function of  $H_\Delta$ , which represents the difference between the actual channel  $H$  and the estimated channel  $\hat{H}$ , i.e.  $H_\Delta = H - \hat{H}$ . It is also assumed that the estimated covariance matrix  $\hat{R}$  is invertible. It also uses the Kullback-Leibler (KL) divergence to measure the distance

between the actual model  $(H, R)$  and the estimated model  $(\hat{H}, \hat{R})$ . Conditioned on the knowledge of  $x$ , the KL divergence between the actual model  $y \sim N(Hx, R)$  and the estimated model  $y \sim N(\hat{H}x, \hat{R})$  takes the form

$$D(f(y|x), \hat{f}(y|X)) = \int \ln \left[ \frac{f(y|x)}{\hat{f}(y|x)} \right] f(y|x) dy \quad (1.21)$$

In this respect, it is worth noting that the KL divergence has been used also in [12] to develop a minimax formulation of robust detection. The use of the KL divergence is rather natural as a metric for model mismatch since it is commonly used by statisticians [13] for fitting statistical models, and by using a differential geometric viewpoint it is argued in [14] that the KL divergence is the natural geometric “distance” between systems.

#### *Nonparametric Uncertainty Models*

It should be emphasized that the classes of allowable characteristics one deals with in robust signal processing are generally nonparametric function classes, such as the class of all power spectral density functions with specified total power (area under the function), which lie between specified upper and lower bounding functions. The uncertainty class for nonparametric uncertainty models can be specified in a variety of ways [15]. It can be based on an  $\epsilon$  – *contamination* model of the type originally proposed by Huber, a total variational model, a spectral band model wherein the power spectral densities specifying the signal and observations are required to stay within a band centered on the nominal Power Spectral Density (PSD) or even, with a probability density model, where the probability density specifying the channel is also required to stay within a constraint set with reference to the nominal Probability Density Function (PDF).

Consider the model

$$Y_k = S_k + N_k. \quad (1.22)$$

The processes  $\{S_k\}$  and  $\{N_k\}$  represent signal and noise, respectively and  $f_S$  and  $f_N$  are the power spectral densities of the signal and the noise. The spectral chosen for designing the estimator may differ from the true signal and noise spectral. This spectral uncertainty is modeled by choosing appropriate classes of PSD’s,  $\mathcal{G}$  and  $\mathcal{N}$ , respectively  $f_S \in \mathcal{G}$  and  $f_N \in \mathcal{N}$ . Next, some specific forms for the uncertainty

classes  $\mathcal{G}$  and  $\mathcal{N}$  are presented. These forms have been widely used to model uncertainty in both the engineering and statistics literature. These forms will be presented for the class  $\mathcal{G}$  but, they could easily be used to model noise spectral uncertainty.

The most commonly used uncertainty class is the  $\epsilon$  – *contaminated* model, which is also known as the  $\epsilon$  – *mixture* or *gross-error* model. This uncertainty class has the following form

$$\mathcal{G}_\epsilon \triangleq \left\{ f_S : f_S(\theta) = (1 - \epsilon)f_S^0(\theta) + \epsilon f'_S(\theta), \quad \forall \theta \in [-\pi, \pi], \right. \\ \left. \int_{-\pi}^{\pi} f'_S(\theta)\lambda(d\theta) = \int_{-\pi}^{\pi} f_S^0(\theta)\lambda(d\theta) \right\} \quad (1.23)$$

where  $f_S^0$  is a nominal PSD and  $\epsilon$  ( $0 \leq \epsilon \leq 1$ ) is the contamination parameter. This class is probably the most popular for representing uncertainty because it models the idea that there is a fraction  $\epsilon$  of completely general uncertainty about the choice of the PSD  $f_S^0$ .

Another common model is the total variational model which has the form

$$\mathcal{G}_{TV} \triangleq \left\{ f_S : \frac{1}{2} \int_{-\pi}^{\pi} |f_S^0(\theta) - f_S(\theta)|\lambda(d\theta) \leq \epsilon, \int_{-\pi}^{\pi} f_S(\theta)\lambda(d\theta) = \int_{-\pi}^{\pi} f_S^0(\theta)\lambda(d\theta) \right\} \quad (1.24)$$

where, again,  $f_S^0$  is a nominal PSD and  $\epsilon$  an uncertainty parameter.

A third model is the band model which has the form

$$\mathcal{G}_B \triangleq \left\{ f_S : f_S^L(\theta) \leq f_S(\theta) \leq f_S^U(\theta), \quad \forall \theta, \int_{-\pi}^{\pi} f_S(\theta)\lambda(d\theta) = 2\pi w \right\} \quad (1.25)$$

where  $\int_{-\pi}^{\pi} f_S^L\lambda(d\theta) \leq 2\pi w \leq \int_{-\pi}^{\pi} f_S^U\lambda(d\theta)$  and  $w$  is the known power of the signal. The name band model comes from the interpretation that  $f_S^L$  and  $f_S^U$  are the lower and upper bounds of a confidence band around a spectral estimate.

A fourth model of interest is the  $p$ -point model which has the form

$$\mathcal{G}_B \triangleq \left\{ f_S : \int_{A_i} f_S(\theta)\lambda(d\theta) = 2\pi w_i, \quad i = 1, \dots, n \right\} \quad (1.26)$$

where the  $A_i$ 's are a partition of  $[-\pi, \pi]$  and  $\sum_{i=1}^n w_i = w$ , the power of the signals. A  $p$ -point class is an appropriate model of uncertainty in situations where, for example, the power  $w_i$  in each interval  $A_i \triangleq [\theta_{i-1}, \theta_i]$  (where  $-\pi = \theta_0 < \theta_1 < \dots < \theta_n = \pi$ ) can be accurately measured using a nested bank of low-pass filters or a bank of bandpass filters.

Note that for each of these classes it is assumed that the power is known. Often it is a reasonable assumption that the power can be accurately estimated even though the shape of the PSD is uncertain.

Nonparametric uncertainty models can also be formulated through probability density functions [8]. Let  $\mathcal{X}$  and  $\mathcal{Y}$  be finite sets denoting the channel input and output alphabets, respectively. The probability law of a (known) channel is specified by a sequence of conditional probability density functions

$$\{Q(y|x) : x \in \mathcal{X}^n, y \in \mathcal{Y}^n\}_{n=1}^{\infty} \quad (1.27)$$

where  $Q(\cdot|\cdot)$  denotes the conditional pdf governing channel used through  $n$  units of time, i.e., “ $n$  uses of the channel.” If the known channel is a discrete memoryless channel, then the law is characterized in terms of a stochastic matrix  $Q : \mathcal{X} \mapsto \mathcal{Y}$  according to

$$Q(y|x) = \prod_{t=1}^n Q(y_t|x_t), \quad (1.28)$$

where  $x = (x_1, \dots, x_n) \in \mathcal{X}^n$  and  $y = (y_1, \dots, y_n) \in \mathcal{Y}^n$ . This kind of uncertainty models assume that the true channel  $Q_{true}(y|x)$  is unknown, though, it is assumed that  $Q_{true}(y|x)$  belongs to a family of channels (1.27), and remains unchanged during the course of a transmission. In [16] a nonparametric uncertainty model is formulated through a joint probability density function. A nominal statistical model is given together with a neighborhood formed by the perturbed models whose KL divergence with respect to the nominal model is bounded by a fixed constant. A static estimation problem is formulated, where the estimate of a random vector  $X \in \mathfrak{R}^n$  needs to be found given an observation model  $Y \in \mathfrak{R}^p$ . It is assumed that the joint nominal density of  $X$  and  $Y$  is Gaussian, so that

$$Z = \begin{bmatrix} X \\ Y \end{bmatrix} \sim N(m_Z, K_Z)$$

where  $m_Z$  and  $K_Z$  denote, respectively, the mean vector and covariance matrix of  $Z$ . Accordingly the nominal probability density of  $Z$  takes the form

$$f_Z(z) = \frac{1}{(2\pi)^{(n+p)/2} |K_Z|^{1/2}} \exp\left(-\frac{1}{2}(z - m_Z)^T K_Z^{-1} (z - m_Z)\right). \quad (1.29)$$

The mean vector and covariance matrices of  $Z$  can be partitioned in terms of the mean vectors and covariances of  $X$  and  $Y$  as

$$m_Z = \begin{bmatrix} m_X \\ m_Y \end{bmatrix}, \quad K_Z = \begin{bmatrix} K_X & K_{XY} \\ K_{YX} & K_Y \end{bmatrix}.$$



$\tilde{f}_Z(z)$  denotes the true probability density of  $Z$ . The KL divergence or relative entropy of  $\tilde{f}_Z$  with respect to  $f_Z$  is given by

$$H(\tilde{f}_Z|f_Z) = \int_{\mathbb{R}^n \times \mathbb{R}^p} \ln \left( \frac{\tilde{f}_Z}{f_Z} \right) \tilde{f}_Z dz. \quad (1.30)$$

The KL divergence is not symmetric and does not obey the triangle inequality, but it satisfies  $H(\tilde{f}_Z|f_Z) \geq 0$  with equality if and only if  $\tilde{f}_Z = f_Z$ .

### 1.1.3 Robust Minimax Estimation

One of the major techniques for designing systems that are robust with respect to modeling uncertainties is the minimax approach, in which the goal is the optimization of worst-case performance. Early applications of game theoretic concepts to communications problems can be found, for example, in the classical estimation studies of Yovits and Jackson [17], and Carlton and Follin [18]. However, the statistical works of Huber in estimation [19] and hypothesis testing [20] are generally regarded as the starting point of the area of minimax robustness, which has been applied successfully to a long sequel of problems in detection and estimation.

There are useful formulations of robustness other than the minimax one, most notably the stability or qualitative robustness ideas introduced by Root [21] in the context of signal detection and by Hampel [22] in the context of parameter estimation. These formulations utilize the idea of robustness as a continuity property of some performance measure as a function of the underlying model. However, from the viewpoint of design, the minimax approach has had the most impact on robust signal processing schemes.

Adaptive procedures, may also be used as robust schemes, when input conditions are not precisely known and may be time-varying. Adaptive procedures, which attempt to learn about input conditions and adjust their specific signal processing structure accordingly to maintain good performance, are generally more complex than fixed minimax schemes. Adaptive schemes are more desirable when the *a priori* uncertainty is so large that the guaranteed level of performance of a minimax scheme would be too poor to be acceptable, and when adequate time and data for adapting are available. Conversely, minimax procedures would be more desirable under more constrained uncertainty classes, and especially as robust procedures

to guard against excessive performance degradation of nominally optimum schemes for deviations from nominal assumptions.

As was remarked, minimax robustness of estimation and hypothesis testing schemes are considered by Huber in [19] and [20], and since then a large number of results on minimax and alternative formulations of robustness have been generated in the statistics literature. In [23] the linear regression problem of estimating an unknown, deterministic parameter vector based on measurements corrupted by colored Gaussian noise is being investigated. Blind minimax estimators are presented and analysed, which consist of a bounded parameter set minimax estimator, whose parameter set is itself estimated from measurements.

The measurements

$$Y = HX + W \quad (1.31)$$

are linear combinations of the parameter vector  $X$ , to which Gaussian noise  $W$  is added. The transformation matrix  $H$  and the noise covariance are assumed to be known. The paper seeks an estimate  $\hat{X}$  which approximates  $X$  in the sense of minimal mean-square error. The parameter to be estimated  $X$  is assumed to lie within a compact parameter set  $S$ . In this case, a linear minimax estimator over the set  $S$  maybe constructed [24], [25]. This is the linear estimator  $\hat{X}_M = GY$  minimizing the worst case MSE among all possible values of  $X$  in  $S$

$$\hat{X}_M = \arg \min_{\hat{X}=GY} \max_{X \in S} E\{\|\hat{X} - X\|^2\}. \quad (1.32)$$

The authors argue that a closed-form solution of (1.32) has been previously derived for many cases of interest and that it has been shown that any linear minimax estimator achieves lower MSE than that of the LS method, for all values of  $X$  in  $S$  [24], [26].

The Blind Minimax Estimators (BMEs) presented in this paper utilize minimax estimators when no parameter set is known. This is done in a two-stage process: 1) A parameter set  $S$  is estimated from the measurements. 2) A minimax estimator designed for  $S$  is used to estimate the parameter vector  $X$ .

The paper investigates two distinct cases. In the first one the authors investigate a spherical blind minimax estimator based on a parameter set of the form  $S = \{X :$

$\|X\|^2 \leq L^2\}$ . For a given value of  $L$  the linear minimax estimator was derived in [25]

$$\hat{X}_M = \frac{L^2}{L^2 + \epsilon_0} \hat{X}_{LS} \quad (1.33)$$

where  $\hat{X}_{LS}$  is the classical LS estimator and  $\epsilon_0$  is the MSE. Now when the bound  $L^2$  is estimated as  $L^2 = \|\hat{X}_{LS}\|^2$ , the Spherical BME (SBME) is given by

$$\hat{x}_{SBME} = \frac{\|\hat{X}_{LS}\|^2}{\|\hat{X}_{LS}\|^2 + \epsilon_0} \hat{X}_{LS}. \quad (1.34)$$

The authors also argue that the BMEs may be constructed around any constant center point  $X_0$ . Furthermore, the paper investigate a second case, the ellipsoidal blind minimax estimator based on a parameter set of the form  $S = \{X : \|X\|_{Q^b}^2 \leq L^2\}$ , for some constant  $b < 0$ . Here  $Q^{-1}$  is the covariance of  $\hat{X}_{LS}$  and the bound  $L^2$  is estimated as  $L^2 = \|\hat{x}_{LS}\|_{Q^b}^2$ . The above cases are being examined in the setting of a linear system of measurements with colored Gaussian noise, and it is shown that the proposed BMEs dominate the LS method, i.e., they achieve lower mean-squared error for any value of the parameter vector. Finally, in [23] the relations of blind minimax techniques to Stein type estimators [27] and least-square regularization is also discussed.

The James-Stein estimator is a nonlinear estimator which can be shown to dominate, or outperform, the “ordinary” (least squares) technique. As such, it is the best-known example of Stein’s phenomenon. In 1961, James and Stein discovered a remarkable estimator that dominates the maximum-likelihood estimate of the mean of a  $p$ -variate normal distribution, provided the dimension  $p$  is greater than two. Various “extended” James-Stein methods were later constructed for the general non-i.i.d. case. However, none of these approaches has become a standard alternative to the LS estimator, and they are rarely used in practice in engineering applications [27].

The problem of estimating an unknown parameter vector  $X$  in a linear model that may be subject to uncertainties, where the vector  $X$  is known to satisfy a weighted norm constraint is also investigated in [25]. This paper addresses two different estimation problems. In the first one, which is similar with the general problem in [23], it is assumed that the model matrix  $H$  is known exactly and the linear estimator that minimizes the worst-case mean-squared error across all possible values of  $X$  is derived. The difference with [23] is that [25] assumes that the parameter set  $S$  is

known, i.e., it is assumed that  $X$  is known to satisfy the weighted norm constraint  $\|X\|_T \leq L$  for some positive definite matrix  $T$  and scalar  $L > 0$ . The robust estimator for this problem is given by (1.33). The second problem, considers the case in which the model matrix  $H$  is subject to uncertainties and seeks the robust linear estimator that minimizes the worst-case MSE across all possible values of  $X$  and all possible values of the model matrix. In many engineering applications, the model matrix  $H$  is subject to uncertainties, for example, it may have been estimated from noisy data, in which case,  $H$  is an approximation to some nominal underlying matrix. If the true data matrix is  $H + \delta H$  for some unknown perturbation matrix  $\delta H$ , then the actual performance of an estimator designed based on  $H$  alone may perform poorly. In this case, robust estimators are considered that explicitly take uncertainties into account. Specifically, in [25] the matrix model  $H$  is not known exactly but is rather given by  $H + \delta H$ , where  $\|\delta H\| \leq p$  and the minimax problem is formulated as

$$\min_{\hat{X}=GY} \max_{\|X\|_T \leq L, \|\delta H\| \leq p} E\{\|\hat{X} - X\|^2\}. \quad (1.35)$$

A similar minimax approach, and model is being used in [28]. Here two problems are being considered. In the first one the uncertainty lies in the covariance of  $X$ ,  $C_X = \tilde{C}_X + \delta C_X$ , where  $\tilde{C}_X$  is known and  $\|\delta C_X\| \leq \epsilon$ . Thus, the minimax MSE problem is formulated as

$$\min_{\hat{X}=GY} \max_{\|\delta C_X\| \leq \epsilon} E(\|\hat{X} - X\|^2). \quad (1.36)$$

In the second problem the covariance matrix  $C_X$  and the model matrix  $H$  are subject to uncertainty. In this case the model matrix  $H$  is given by  $H = \tilde{H} + \delta H$ ,  $\|\delta H\| \leq p$ , where  $\tilde{H}$  is known. The minimax problem is formulated as

$$\min_{\hat{X}=GY} \max_{\|\delta C_X\| \leq \epsilon, \|\delta H\| \leq p} E(\|\hat{X} - X\|^2). \quad (1.37)$$

In order to improve the performance over the minimax MSE approach, a competitive approach is also considered in [28]. This approach assumes that  $H$  is completely known, and seeks a linear estimator  $\hat{X}$  that minimizes the worst case regret in order to partially compensate for the conservative character of the minimax approach. The regret  $\mathcal{R}(C_X, G)$  is defined as the difference between the MSE using an estimator  $\hat{X} = GY$  and the smallest possible MSE attainable with an estimator of the form  $\hat{X} = G(C_X)Y$  when the covariance  $C_X$  is known, which is denoted as  $MSE^0$ . The

minimax problem is formulated as

$$\min_G \max_{\| \delta C_X \| \leq \epsilon} \mathcal{R}(C_X, G) \quad (1.38)$$

where  $\mathcal{R}(C_X, G) = E(\|\hat{X} - X\|^2) - MSE^0$ . The linear minimax regret estimator is shown to be equal to a Minimum-Mean-Squared-Error (MMSE) estimator corresponding to a certain choice of signal covariance, that depends explicitly on the uncertainty region.

A different approach is introduced in [16], where given a nominal statistical model, the minimax estimation problem consisting of finding the best least-squares estimator for the least favorable statistical model within a neighborhood of the nominal model is being considered. The uncertainty model for this problem has already been described in Section 1.1.2. The neighborhood is formed by placing a bound on the KL divergence between the actual and nominal models described by the uncertainty set  $\mathcal{B} = \{\tilde{f}_Z \in \mathcal{F} : H(\tilde{f}_Z | f_Z) \leq c\}$ .

The authors seek to estimate a random vector  $X \in \mathfrak{R}^n$  given an observation vector  $Y \in \mathfrak{R}^p$  such that the joint nominal density of  $Z = [X^T, Y^T]^T$  is Gaussian with the parameterization (1.29). They need to find an estimator  $\hat{X} = g(Y)$  and a least favorable density  $\tilde{f}_Z$  which solve the minimax problem

$$\min_{g \in \mathcal{G}} \max_{\tilde{f}_Z \in \mathcal{B}} J(\tilde{f}_Z, g) \quad (1.39)$$

where

$$J(\tilde{f}_Z, g) = \frac{1}{2} E_{\tilde{f}_Z} [\|X - g(Y)\|^2] = \frac{1}{2} \int \|x - g(y)\|^2 \tilde{f}_Z(z) dz \quad (1.40)$$

For the above minimax problem it is shown that  $J(\tilde{f}_Z, g)$  admits a saddle point and that the estimator takes the form

$$g_0(Y) = m_X + G_0(Y - m_Y) \quad (1.41)$$

with  $G_0 = K_{XY} K_Y^{-1}$ . This estimator is linear and yields the usual least-squares estimate of  $X$  given  $Y$  for both the nominal density  $f_Z$  and the perturbed density  $\tilde{f}_0$ . In [16] is shown that for a Gaussian nominal model and a finite observations interval, or for a stationary Gaussian process over an infinite interval, the usual non-causal Wiener filter remains optimal. However, it also shows that in the causal case, the usual least-squares estimator or the causal Wiener filter are no longer optimal, and

a characterization is given for the structure of an optimal robust estimator and the matching of least favorable statistical model. The optimal causal filter derived is a risk-sensitive filter, where the risk-sensitivity parameter was selected to match the maximum allowable relative entropy for the perturbed model.

The same minimax approach used in [16] is also applied in [11], the difference is in the uncertainty model used. As described on Section 1.1.2, [11] uses a parametric uncertainty model, where the uncertainty is placed in the channel  $H$  and noise  $R$  information. Similar with [16], the KL divergence is used to measure the “distance” between the actual and the estimated model, and the uncertainty set is defined as  $\mathcal{B} = \{(H_\Delta, R) : D(H, R; \hat{H}, \hat{R}) \leq c\}$ . Then the MMSE equalizer problem is formulated by finding as

$$\min_F \max_{(H_\Delta, R) \in \mathcal{B}} J(F, H_\Delta, R) \quad (1.42)$$

where  $F$  is an MMSE equalizer matrix, and  $(F, H_\Delta, R)$  is the MSE pay-off function. By using Lagrangian duality, the minimax problem is transformed into an equivalent min-min problem over a convex domain, where the standard convex optimization methods apply. Then, it is shown that the robust MSE equalizer can be obtained by solving numerically a scalar convex minimization problem.

Although robustness issues do not appear explicitly in [29], and the minimax approach is not really implemented per say, it presents some interesting results for filtering problems when there is uncertainty in the exact value of the probability model. In the problems described above the minimax approach first maximizes the MSE in order to derive the worst case measure given a specific uncertainty set. In [29] given a measurable space  $(\Omega, \mathcal{F})$  and random variables  $X, Y$ , a function of  $X$ ,  $\phi = \phi(X)$  needs to be estimated by random variable  $\hat{\phi} \in \mathcal{Y}$ , where  $Y$  represents the observations or measurements. The true distribution  $P_{\alpha_0}$  is assumed to belong to a family of probability measures  $\{P_\alpha\}_{\alpha \in A}$ . The minimum cost estimator is defined as the estimator which minimizes the error cost function  $E_{\alpha_d}[p_\alpha(\phi - \hat{\phi})]$ , where  $p_\alpha$  is a strictly convex function and  $P_{\alpha_d}$  is assumed as the true model. Instead of maximizing the above error cost function, it chooses an exponential cost function. As stated in [29], it is known from control theory that if a controller is designed to minimize an average-of-exponential criterion then this controller implies certain robustness properties. Therefore, a minimum risk-sensitive estimator, denoted  $\hat{\Phi}_{rs}^*$ , is defined as the minimum cost estimator obtained by selecting  $p_\alpha(e) = \exp(\mu p(e))$ , where  $\mu > 0$

is a parameter determining the degree of “risk”,

$$\hat{\Phi}_{rs}^* = \arg \min_{\hat{\phi} \in R} \int \exp(\mu p(\phi - r)) \pi_{\alpha_d}(dx) \quad (1.43)$$

where  $\pi_{\alpha_d}$  is the conditional distribution of  $X$  given  $Y$  under  $P_{\alpha_d}$ .

Furthermore, it is shown that the risk-sensitive estimators enjoy an error bound which is the sum of two terms, the first of which coincides with an upper bound on the error one would obtain if one knew exactly the underlying probability model, while the second term is a measure of the distance between the true and design probability models. The first term quantifies “good performance” under nominal conditions, and the second term quantifies the “acceptable performance” under non-nominal conditions. Also, the second term plays a major role in determining the class of permissible variations from nominal.

Next, [30] deals with the problem of designing robust linear causal estimators of linear functions of discrete-time wide-sense stationary signals, when the knowledge of the signal and/or noise spectral is inexact. The spectral uncertainty is modeled by choosing appropriate classes of PSD’s,  $\mathcal{G}$  and  $\mathcal{N}$ , and assuming that the signal PSD  $f_S \in \mathcal{G}$  and the noise PSD  $f_N \in \mathcal{N}$ . The paper seeks to find the transfer function  $H_R^*$  with the smallest possible upper bound on the MSE,  $e_D(f_S, f_N; H_R)$  over all  $f_S$  in  $\mathcal{G}$  and  $f_N$  in  $\mathcal{N}$ , that is the solution to the following minimax formulation

$$\inf_{H \in H_+^2} \sup_{(f_S, f_N) \in \mathcal{G} \times \mathcal{N}} e_D(f_S, f_N; H) \quad (1.44)$$

The solution is given under mild regularity conditions in the terms of the least favorable pair of spectra,  $(f_S^L, f_N^L)$ , thus reducing the minimax problem to a direct maximization problem which in many cases can be solved easily. The solution is based on the fact that a pair of PSD’s  $(f_S^L, f_N^L) \in \mathcal{G} \times \mathcal{N}$  and its optimal causal transfer function  $H_L^*$ , under specific conditions, form a saddle-point solution to (1.44) if and only if  $(f_S^L, f_N^L)$  is least favorable for causal estimation. Furthermore, solutions are given explicitly for the problem of robust causal filtering of an uncertain signal in white noise.

#### 1.1.4 Relation of Mutual Information and MMSE

As was mentioned in Section 1.1.1 one of the most popular Bayesian estimation techniques is the Least-Square estimation, also known as Minimum-Mean-Squared-

Error (MMSE) estimation. Recently, extensive work has been done which connects the mutual information between the input and the output of a channel, and the MMSE in estimating the input given the output. In a wider context, the mutual information and mean-square error are at the core of information theory and estimation theory, respectively. The input-output mutual information is an indicator of how much coded information can be pumped through a channel reliably given a certain input signaling, whereas the MMSE measures how accurately each individual input sample can be recovered using the channel output.

The relation between the mutual estimation and MMSE is presented in [31], which deals with arbitrarily distributed finite power input signals observed through an additive Gaussian noise channel. A new formula is presented that connects the input-output mutual information and the minimum mean-square error achievable by optimal estimation of the input given the output. Given that the input-output mutual information and the MMSE are monotone functions of the signal-to-noise ratio (SNR), denoted by  $I(snr)$  and  $mmse(snr)$ , respectively, the mutual information in nats and the MMSE satisfy the following relationship regardless of the input statistics:

$$\frac{d}{dsnr} I(snr) = \frac{1}{2} mmse(snr) \quad (1.45)$$

That is, the derivative of the mutual information (nats) with respect to the signal-to-noise ratio (SNR) is equal to half the MMSE, regardless of the input statistics. This relationship holds for both scalar and vector signals.

For the scalar signal, [31] considers a pair of real-valued random variables related by

$$Y = \sqrt{snr}X + N \quad (1.46)$$

where  $snr \geq 0$  and  $N \sim N(0, 1)$  is a standard Gaussian random variable independent of  $X$ . Then  $X$  and  $Y$  can be regarded as the input and output, respectively, of a single use of a scalar Gaussian channel with an SNR of  $snr$ .

On the other hand, the multiple-input multiple-output (MIMO) system is described in [31] by the vector Gaussian channel

$$Y = \sqrt{snr}HX + N \quad (1.47)$$



where  $H$  is a deterministic  $L \times K$  matrix and the noise  $N$  consists of independent standard Gaussian entries. The input  $X$  (with distribution  $P_X$ ) and the output  $Y$  are column vectors of appropriate dimensions.

Various applications of the above relationship are identified in [31], e.g. in relating code-division multiple-access (CDMA) channel spectral efficiencies (mutual information per dimension) under joint and separate decoding in the large-system limit. Also the fact that the mutual information and the MMSE determine each other by a simple formula provides new means to calculate or bound one quantity using the other. An upper bound for the mutual information is immediate by bounding the MMSE for all SNRs using a suboptimal estimator. Lower bounds on the MMSE lead to new lower bounds on the mutual information.

The low-SNR asymptotics of  $mmse(X; snr)$  is also studied extensively in [31], mainly for the scalar Gaussian channel. The Taylor expansion of  $mmse(X; snr)$  at  $snr = 0$  is obtained and the coefficients turn out to depend only on the moments of  $X$ . Based on the above work [32] deals with the high-SNR asymptotics of MMSE in Gaussian channels and defines a new information measure called MMSE dimension. The MMSE dimension of  $X$  for the scalar channel (1.46) is defined as the limit as  $snr \rightarrow \infty$  of the product of  $snr$  and the MMSE. For discrete, absolutely continuous or mixed  $X$  [32] shows that the MMSE dimension equals Rényi's information dimension. For singular  $X$ , it shows that the product of  $snr$  and MMSE oscillates around information dimension periodically in  $snr$  (dB).

Finally, based on the relationship (1.45), [33] presents some new results on mutual information and various regularity properties of the MMSE functional are explored together with its connections to Shannon theory.

## 1.2 Thesis Motivation

Robust estimation is not a new subject of study, as was already seen, and it is a major problem in statistics and signal processing. The minimax approach, in which the goal is to optimize the worst case performance, is one of the major techniques for designing robust systems with respect to modeling uncertainties and has been applied to many problems in detection and estimation, as has already presented in

Section 1.1.3. This thesis will focus, mainly, on least-square estimation problems when the statistics or internal dynamics describing the signal and observations are not known exactly. A robust estimation approach is being proposed in this thesis which employs a minimax formulation.

A measurable space  $(\Omega, \mathcal{F})$  is given on which the unobserved Random Variable,  $X$  and the observed RV  $Y$  are defined, via  $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \Sigma_{\mathcal{X}})$ ,  $Y : (\Omega, \mathcal{F}) \rightarrow (\mathcal{Y}, \Sigma_{\mathcal{Y}})$ . The objective is to estimate  $X$  by a function of the random variable  $Y$ . As mentioned in Section 1.1.1, the classical least square estimation problem deals with minimization of the average pay-off which can be expressed in the following ways

$$J(\Phi^*) = \inf_{\Phi \in \mathcal{X}_{ad}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) dP_{X,Y}(x, y) \quad (1.48)$$

$$= \inf_{\Phi \in \mathcal{X}_{ad}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \mu(x, dy) dP_X(x) \quad (1.49)$$

$$= \inf_{\Phi \in \mathcal{X}_{ad}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \eta(y, dx) dP_Y(y). \quad (1.50)$$

Clearly if  $\ell(x, \Phi(y)) = \|x - \Phi(y)\|_{\mathbb{R}^n}$ , then the estimate of  $X$  denoted as  $\hat{X} = \Phi(Y)$  is given by  $\Phi^*(Y) = E[X|Y]$ . As was already mentioned, and shown from the literature the signal statistics  $P_{X,Y}$  or  $\mu, \eta$  are not always known. In this thesis the uncertainty description of the system, and the nominal description of the system are modeled by probability distributions, or general measures, defined on measurable spaces. Moreover, two type of uncertainty models are being considered.

1. Uncertainty Models on Conditional Distributions or otherwise known Stochastic Kernels;
  - i) When the conditional probability distribution of the measurement  $Y$  given the signal to be estimated  $X$ , or channel kernel, is unknown;
  - ii) When the *a posteriori* distribution of  $X$  given  $Y$  is unknown;
2. Uncertainty Models on Joint Distributions.

The goal is to derive new robust estimators for the above uncertainty models (this is described explicitly in Chapter 3). Notice that (1.48) will be used when the uncertainty is on the joint distribution, while (1.49), (1.50) will be used when the uncertainty is on the channel kernel, *a posteriori* distribution, respectively.

Stochastic kernel uncertainty models are appropriate for communication system design, in which the input message has a known distribution, while the channel is unknown but belongs to a certain class of channels. These are nonparametric uncertainty models which so far have not been taken into consideration. Joint distribution uncertainty models are usually employed when both the unobserved and observed random variables are uncertain. Joint uncertainty models are used in [16] and [29], but [29] does not implement a minimax approach and it shows mainly that the mean-square error for the true model is bounded by the sum of two terms, with the first term representing the performance of the risk-sensitive filter with respect to the nominal model, and the second term corresponding to the relative entropy between the actual and nominal models. On the other hand [16] employs a minimax viewpoint but uses the additional assumption that the system considered is Gaussian and it uses a parametric approach to derive its results.

The uncertainty description of the above systems is characterized by the class of uncertain measures which satisfy a relative entropy constraint with respect to a nominal measure. The use of the relative entropy, also known as the KL divergence, is rather natural as a metric for model mismatch since it is commonly used by statisticians for fitting statistical models, and by using a differential geometric viewpoint it is argued that the KL divergence is the natural geometric “distance” between systems. The KL distance constraint has been also used in robust estimation problems in [11], [16] and [29].

By using the KL divergence various constraint sets are defined for each uncertainty model, and for each the minimax estimation problem is solved and the worst case measure of the true model and also the robust estimator are derived. This is done in a generalized framework, while the estimators derived can be used accordingly to various estimation applications.

One particular application which is investigated is the robust nonlinear estimation for finite-dimensional autoregressive channel models found in [34], [35] subject to an uncertainty set. Autoregressive channel models have been used with success to predict fading channel dynamics for the purposes of Kalman filter based channel estimation and for long-range channel forecasting. They have also been used by several authors to simulate correlated Rayleigh fading. A specific example is the non-coherent estimation problem, where there is an attenuated sinusoid in a multipath environment which is subject to an additive Gaussian noise. Here an uncertainty

model on the joint distribution is used and as before, the uncertainty set is being described by a KL constraint.

The above problem, which is a nonlinear problem, will be addressed through a minimax approach. The minimax approach has not been used so far for this specific problem and the difficulty lies in the minimization, and the way of obtaining the solution of the  $\hat{a}$  posteriori distribution recursion in closed form. In general when dealing with nonlinear estimation problems, approximations have to be made to find sub-optimal nonlinear estimators. The standard method is to use the Taylor series expansion and apply linear filtering theory, giving rise to the so-called extended Kalman filter (EKF) [5, 6, 7]. Other more sophisticated sub-optimal estimation techniques are available, e.g. reiteration, higher order filters, and statistical linearization [5]. The scope of the thesis is not to investigate general nonlinear estimation techniques but to find a way to tackle the minimization in the minimax formulation. Here the maximization is addressed using variational methods, while the minimization is addressed using a change of probability measure technique [36]. The change of probability measure techniques introduced is being used in order to derive recursive equations for the conditional distribution of nonlinear filtering problems, which helps us to compute the worst-case pay-off functional. Minimizing this worst-case pay-off function gives the desired robust minimax estimators. Special emphasis is given also to the connection between the robust minimax estimation problem and the classical non-coherent estimation.

Even though the primal focus of this thesis is robust least square estimation some other well-known estimation techniques, like the MAP and the ML estimation techniques, are also investigated. As was described in Section 1.1.1, in the classical MAP problem, by applying a uniform cost criterion, the estimator of  $x$  is derived through the maximization of the  $\hat{a}$  posteriori (discrete or continuous) density  $\eta(y, dx)$ . In this thesis instead of a uniform cost criterion, an exponential one is being used. As its been mentioned in [29], it is known in control theory that if a controller is designed to minimize an average-of-exponential (or risk-sensitive) criterion, then this controller might have certain robustness properties. Hence, it is natural to consider the use of such exponential criteria in filter design. The objective here is to derive a generalized MAP estimator using an exponential cost criterion and also show a connection with the minimax approach used in the least-square estimation problem. A similar approach is also used for the derivation of generalized ML estimator. The ML estimation technique, which is not a Bayesian approach, is an alternative method for

deriving estimators, when the parameter  $x$  is assumed to be unknown but without being characterized by a probabilistic structure. The classical method of maximum likelihood estimates  $x$  by finding the value of  $x$  that maximizes the likelihood function  $\mathcal{L}(x) = \mu(x, dy)$ . The goal is to derive a generalized ML estimator when this likelihood function includes also an exponential function and show the connection with the generalized MAP estimator.

### 1.3 Thesis Objective

The main objective of this thesis is to derive robust estimators for least-square estimation problems in the presence of model uncertainties. Additional to this objective is the investigation of a generalized MAP and ML estimation technique. These three estimation techniques are well known and often used in signal estimation and in signal processing.

Firstly, least-square estimation problems are investigated when the real description of the system is unknown and the only knowledge is that it belongs to an uncertainty set, or a class of systems. The real description of the system, the uncertainty description of the system, and the nominal description of the system are modeled by probability distributions, or general measures, defined on measurable spaces. The uncertainty is described by a KL constraint between the unknown distribution and a fixed nominal distribution.

Secondly, the behavior of the Maximum  $\hat{A}$  Posteriori estimation and the Maximum Likelihood estimation technique is investigated, when introducing new elements in the derivation of the classical estimators. Generalized estimators for both techniques are derived by altering the cost function for the MAP case and the likelihood function for the ML case.

Chapter 2 presents background material and explains the main techniques used throughout this thesis, which are change of probability theory and the minimax theory.

Chapter 3 considers least-square estimation problems for classes of models and introduces the concept of uncertainty. The uncertainty is described by a KL distance

constraint between the unknown distribution and a fixed nominal distribution. The theory and contribution of this chapter are developed at two levels of generality; the abstract level and the application level. At the abstract level the uncertainly models used are Stochastic Kernels and Joint Distributions.

Starting with the first uncertainty model, the relation between the unobserved RV  $X$  and the observed RV  $Y$  is defined via a probabilistic mapping,  $\mu : \mathcal{X} \times \Sigma_Y \rightarrow [0, 1]$ , which satisfies the following two conditions:

1. For every  $x \in \mathcal{X}$ , the set function  $\mu(x, \cdot)$  is a probability measure on  $\Sigma_Y$  (possibly finite additive);
2. For every  $F \in \Sigma_Y$ , the function  $\mu(\cdot, F)$  is  $\mathcal{X}$ -measurable.

The mapping  $\mu$  is called a stochastic kernel or transition probability and represents the nominal system model or mapping, which is fixed. The true kernel denoted by  $\nu : \mathcal{X} \times \Sigma_Y \rightarrow [0, 1]$  is assumed unknown. Envisioned scenarios are communication channels whose nominal behavior is known, while its true conditional distribution is unknown. The KL distance is used as a measure of distance between the true model and uncertainty model, hence the true kernel is assumed to belong to the pointwise uncertainty set,

$$\mathcal{A}^x(\mu) \triangleq \left\{ \nu \in \mathcal{P} : H(\nu|\mu)(x) \leq R(x) \right\} \quad (1.51)$$

where  $R : \mathcal{X} \rightarrow [0, \infty)$  and  $H(\cdot|\cdot)(x) : \mathcal{X} \rightarrow [0, \infty]$  is the KL distance between two kernels. Additionally, the following uncertainty set is defined

$$\mathcal{A}(\mu) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{X}} H(\nu|\mu)(x) dP_X(x) \leq \int_{\mathcal{X}} R(x) dP_X(x) \stackrel{\nabla}{=} r_1 \right\}. \quad (1.52)$$

For the second uncertainty model, uncertainty on the  $\hat{a}$  posteriori distribution, the mapping  $\eta : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$  is considered, that satisfies the following two conditions:

1. For every  $y \in \mathcal{Y}$ , the set function  $\eta(y, \cdot)$  is a probability measure on  $\Sigma_{\mathcal{X}}$ ;
2. For every  $F \in \Sigma_{\mathcal{X}}$ , the function  $\eta(\cdot, F)$  is  $\mathcal{Y}$ -measurable.

This probabilistic kernel  $\eta(y, dx)$  represents the nominal system model or mapping ( $\hat{a}$  posteriori information) and the true kernel  $\nu(y, dx)$ , denoted by  $\nu : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$

belongs to an uncertainty set. The following two uncertainty sets are considered.

$$\mathcal{B}^y(\eta) \triangleq \left\{ \nu \in \mathcal{P} : H(\nu|\eta)(y) \leq R(y) \right\} \quad (1.53)$$

$$\mathcal{B}(\eta) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{Y}} H(\nu|\eta)(y) dP_Y(y) \leq \int_{\mathcal{Y}} R(y) dP_Y(y) \stackrel{\nabla}{=} r_2 \right\} \quad (1.54)$$

where  $R : \mathcal{Y} \rightarrow [0, \infty)$  and  $H(\cdot|\cdot)(y) : \mathcal{Y} \rightarrow [0, \infty]$  is the KL distance between two kernels. Additionally, the next case is investigated, where the true kernel  $\nu(y, dx)$  belongs to a new uncertainty set described by

$$\mathcal{B}^R(\eta) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{Y}} H(\nu|\eta)(y) dP_Y(y) \leq \int_{\mathcal{X} \times \mathcal{Y}} R(x, y) \nu(y, dx) dP_Y(y) + \bar{R} \stackrel{\nabla}{=} r_3 \right\} \quad (1.55)$$

where  $R : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty)$ ,  $R \in BC(\mathcal{X} \times \mathcal{Y})$  and  $H(\cdot|\cdot)$  is the KL distance between two kernels.

Finally, uncertainty is modeled via the joint distribution of  $X$  and  $Y$ . This model is appropriate when one wishes to model *a priori* uncertainty. It is assumed that the joint distribution  $P_{X,Y}(x, y)$  represents the nominal system model and the true joint distribution denoted by  $Q \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ , is assumed to belong to an uncertainty set described by

$$\mathcal{C}(P_{X,Y}) \triangleq \{Q_{X,Y} \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y}) : H(Q_{X,Y}|P_{X,Y}) \leq R\}$$

where  $R \in [0, \infty)$  and  $H(\cdot|\cdot)$  is the KL distance between the two joint distributions.

For all the above uncertainty sets the goal is to formulate the minimax problem and derive first the worst-case measure for the true stochastic kernel  $\nu^*$  or true joint distribution  $Q^*$ , and then to derive the robust estimator for each case. The minimax estimation problem can be formulated as

$$\begin{aligned} J_1(\Phi, \nu) &\triangleq \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x), \\ J_1(\Phi^*, \nu^*) &= \inf_{\Phi \in \mathcal{X}_{ad}} \sup_{\nu \in \mathcal{A}(\mu)} J_1(\Phi, \nu). \end{aligned} \quad (1.56)$$

This represents one of the uncertainty models. Similar formulation will be used for all the other cases.

First, the appropriate space of measures is introduced and then the maximizing kernel and joint measure are computed explicitly using Lagrangian functionals, and

variational methods. Moreover, important monotonicity properties satisfied by the optimal strategies are presented, which can be used to develop numerical algorithms for computation of the optimal solution, and upper and lower bounds on the optimal solution. Furthermore, in Chapter 3, several problems of estimation theory are formulated, and solutions are sought when the models (conditional distributions, joint distributions) are uncertain, and they belong to specific subsets of the set of conditional or joint distributions. These examples include MIMO communication systems.

Chapter 4 applies the theory developed in Chapter 3 to finite-dimensional autoregressive channel models, and derives robust estimators for a class of uncertain models. Similar to Chapter 3, the uncertainty is described by a relative KL distance between the unknown joint distribution and a fixed nominal joint distribution. The abstract channel model is given by

$$\begin{aligned} x_{k+1} &= f(k+1, x_k) + B_{k+1}w_{k+1}, \quad x_0 \in \mathfrak{R}^n \\ y_k &= h(k, x_k) + D_k v_k, \quad y_0 \in \mathfrak{R}^d. \end{aligned} \quad (1.57)$$

Here  $x_0 : \Omega \rightarrow \mathfrak{R}^n$  is the initial state and  $w : \Omega \times N_0 \rightarrow \mathfrak{R}^n$ ,  $v : \Omega \times N_0 \rightarrow \mathfrak{R}^d$ , are random noises, all mutually independent.

The uncertainty model being used here is the joint distribution, where  $P_{x^m, y^m}$  denotes the nominal (in the absence of modeling uncertainties) joint distribution of the sequences  $(x^m, y^m)$ , which corresponds to the one induced by model (1.57) and  $Q_{x^m, y^m}$  denotes the true joint distribution of the sequences  $(x^m, y^m)$ , which is unknown. The only available information is that  $Q_{x^m, y^m}$  belongs to a class of possible distributions. This class is modelled by the information theoretic KL distance set

$$\mathcal{C}(P_{x^m, y^m}) \triangleq \{Q_{x^m, y^m} : H(Q_{x^m, y^m} | P_{x^m, y^m}) \leq R\}.$$

The minimax estimation problem is defined as

$$J(\tilde{x}^*, Q_{x^m, y^m}^*) = \inf_{\tilde{x}^m \in \mathcal{X}_{ad}} \sup_{Q_{x^m, y^m} \in \mathcal{C}(P_{x^m, y^m})} E_{Q_{x^m, y^m}} \left\{ \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right\} \quad (1.58)$$

where  $\tilde{\ell}(x_k, \tilde{x}_k)$  is a measure of distance between the state  $x_k$  and its estimate  $\tilde{x}_k$ .

In the abstract setting, the maximization is addressed using variational methods, while the minimization is addressed using a change of probability measure technique



[36]. Through this technique conditional expectations are related via

$$E \left[ \Phi(x_m) \exp \left( \frac{1}{s} \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right) \middle| \mathcal{Y}_m \right] = \frac{\bar{E} \left[ \Phi(x_m) \Lambda_m \exp \left( \frac{1}{s} \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right) \middle| \mathcal{Y}_m \right]}{\bar{E} \left[ \Lambda_m \middle| \mathcal{Y}_m \right]} \quad (1.59)$$

where  $\Lambda$  is the likelihood function of the complete data and  $E, \bar{E}$  denotes expectation under the probability distribution  $P, \bar{P}$ , respectively. The change of probability measure techniques introduced is being used in order to derive recursive equations for the conditional density of nonlinear filtering problems.

Next, the theory developed is applied to two specific applications. Firstly to a linear Gaussian model and secondly to an attenuated sinusoid in a multipath environment which is subject to an additive Gaussian noise given by

$$y(t_k) = \sum_{i=1}^N A_i(t_k) r_i \cos(\omega_c(t_k - \tau_i(t_k)) + \theta_i) S(t_k - \tau_i(t_k)) + D(t_k) v(t_k). \quad (1.60)$$

Various estimators are being derived for both applications. Furthermore, for the attenuated sinusoid model a connection to least-squares estimation found in [1], [4] for single channel is established by reducing the uncertainty to zero, while generalizing existing results.

Chapter 5 considers the classical Maximum  $\hat{A}$  Posteriori estimation and Maximum Likelihood estimation problems. The classical MAP estimator is derived by minimizing a specific uniform cost function. The theory behind maximum  $\hat{a}$  posteriori estimation is being investigated when this uniform cost function is modified. An exponential cost function is introduced defined by

$$C(X, \Phi(y)) = \begin{cases} e^{\frac{\ell(x)}{s}}, & \text{if } |X - \Phi(y)| > \Delta \\ 0, & \text{if } |X - \Phi(y)| \leq \Delta \end{cases}$$

where  $\ell : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$  is an  $\Sigma_{\mathcal{X}} \times \Sigma_{\mathcal{X}}$ -measurable function and  $\Delta > 0$ . By minimizing the average cost function a generalized MAP estimator is derived. Furthermore, a connection to the theory of minimax least-square estimation derived in Chapter 3 is presented.

The classical ML estimator is derived by maximizing the maximum likelihood function. The theory behind ML estimation is investigated, when a modified likelihood

function is being used, defined as  $e^{\frac{\ell(x)}{s}} \mu(x, dy)$ , and a generalized ML estimator is derived. A relation with the MAP estimation is also presented. Finally, the theory developed is applied to several examples.

Chapter 6 contains the main points of the thesis and suggests directions for future research.

## 1.4 Contributions

The main contributions of Chapter 3 are the following.

1. New uncertainty models based on stochastic kernels and joint distributions are presented for robust minimax least-square estimation problems;
2. For each uncertainty model, various uncertainty sets are defined using the KL distance constraint;
3. The maximization problem is addressed using variational methods and for each uncertainty model the worst case measure is derived;
4. The worst case pay-off is derived, for each uncertainty model, and by minimizing this pay-off the robust estimators are derived;
5. The derived worst case measures and worst case pay-offs are employed to compute robust estimators for linear problems, which also include examples from MIMO communication systems.

The main contributions of Chapter 4 are the following:

1. A minimax least-square estimation problem is formulated for finite-dimensional autoregressive channel models, when the uncertainty inserted is on the joint distribution, and the uncertainty set is described by the KL distance;
2. The minimization is addressed using a change of probability measure technique, and a recursive equation for the conditional distribution is derived. Using this unnormalized  $\hat{a}$  posteriori distribution, the worst case pay-off is derived;

3. The derived unnormalized  $\hat{a}$  posteriori distribution is employed to compute the estimator for a linear gaussian autoregressive channel;
4. The derived unnormalized  $\hat{a}$  posteriori distribution is employed to compute robust phase and envelope estimators for an attenuated sinusoid in a multipath environment, which is subject to an additive Gaussian noise;
5. A connection between the derived robust results and the classical non-coherent estimation problem is presented.

The main contributions of Chapter 5 are the following:

1. The cost function used in the MAP estimation method is defined by an exponential function;
2. A generalized MAP estimator is derived using the exponential cost function;
3. A generalized ML estimator is derived by inserting an exponential function in the maximum likelihood function;
4. The derived generalized estimators are used in several examples.



# CHAPTER 2

## BACKGROUND MATERIAL

In this chapter, background material is presented and the main techniques used throughout this dissertation are explained. In Section 2.1 the basic mathematical concepts of vector spaces, probability theory, are reviewed following a measure-theoretic approach and the theory behind minimax estimation techniques is described. In Section 2.2, the theory of change of probability measure technique is presented.

### 2.1 Basic Mathematical Background

In the study of systems, functional analysis plays a fundamental role. The concepts and mathematical tools developed in this section are used in subsequent chapters.

#### 2.1.1 Functional Analysis

Here, the basic elements of mathematics that are needed to understand systems that can be modeled by linear or nonlinear differential equations are presented. To fully describe the state of the process of a system at any point of time a number of variables have to be quantified which are called a vector. Moreover, some important results from functional analysis are presented, which will be frequently used in the following chapters.

**Definition 2.1.1.** (*Vector Space*) A vector space over a field  $F$  is a set of vectors  $\mathcal{Z}$  together with the operations of

- addition in  $\mathcal{Z}$ :  $x + y$ ,  $x, y \in \mathcal{Z}$ ,
- multiplication by Scalar:  $\alpha \cdot x$ ,  $\alpha \in F$ ,  $x \in \mathcal{Z}$ ,

satisfying the following properties:

1. *Associativity*:  $(x + y) + z = x + (y + z)$ ,  $\forall x, y \in \mathcal{Z}$ .
2. *Commutativity*:  $x + y = y + x$ ,  $\forall x, y \in \mathcal{Z}$ .
3. *Additive identity*: There exists  $\emptyset \in \mathcal{Z}$  such that  $\emptyset + x = x$ ,  $\forall x \in \mathcal{Z}$ .
4. *Additive inverse*:  $\forall x \in \mathcal{Z}$  there exist  $-x \in \mathcal{Z}$ , such that  $x + (-x) = \emptyset$ ,  $\forall x \in \mathcal{Z}$ .
5. *Ass. Scalar*:  $\alpha(\beta \cdot x) = (\alpha \cdot \beta)x$ ,  $\forall x \in \mathcal{Z}$ ,  $\alpha, \beta \in F$ .
6. *Multiplicative identity*:  $1 \cdot x = x$ ,  $0 \cdot x = \emptyset$ ,  $\forall x \in \mathcal{Z}$ .
7. *Scalar Mult. Distributive w.r.t. vector addition*:  $\alpha(x + y) = \alpha x + \beta y$ ,  $\forall x, y \in \mathcal{Z}$ .
8. *Scalar Mult. is Distr. w.r.t. scalar addition*:  $(\alpha + \beta)x = \alpha x + \beta x$ ,  $\forall \alpha, \beta \in F$ .

Usually, a vector space over  $F = \mathbb{R}$  is called a real vector space and a vector space over  $F = \mathbb{C}$  is called a complex vector space.

**Definition 2.1.2.** (Normed Space) A Normed Space is a vector space  $\mathcal{Z}$  furnished with a norm  $\|\cdot\|_{\mathcal{Z}}$  and denoted by  $(\mathcal{Z}, \|\cdot\|_{\mathcal{Z}})$ . The norm  $\|\cdot\|_{\mathcal{Z}} : \mathcal{Z} \rightarrow \mathbb{R}$  is a real valued function defined on  $\mathcal{Z}$  which must satisfy the following properties [37]:

- i)  $\|x\| \geq 0 \forall x \in \mathcal{Z}$ ;
- ii)  $\|x\| = 0$  if and only if  $x = 0$ ;
- iii)  $\|\alpha x\| = |\alpha| \|x\|$ ,  $\forall x \in \mathcal{Z}$ ,  $\alpha \in F$ ;
- iv)  $\|x + y\| \leq \|x\| + \|y\|$ ,  $\forall x, y \in \mathcal{Z}$ .

**Definition 2.1.3.** (Cauchy sequence) Let  $(\mathcal{Z}, \|\cdot\|_{\mathcal{Z}})$  be a normed space. A sequence  $\{x_n\} \in \mathcal{Z}$  is said to be a Cauchy sequence if [37]

$$\lim_{n \rightarrow \infty} \|x_{n+p} - x_n\| = 0 \text{ for every } p \geq 1.$$

**Definition 2.1.4.** (Banach Space) A normed space  $\mathcal{Z}$  is said to be complete if every Cauchy sequence of  $\mathcal{Z}$  has a limit on  $\mathcal{Z}$ . A complete normed space is called a Banach space [37].

**Definition 2.1.5.** (Hilbert Space) A Hilbert space  $\mathcal{H}$ , is a Banach space with special structure. It is furnished with an inner (or scalar) product  $(\cdot, \cdot)$  as defined below [37].

For any two elements  $f, g \in \mathcal{H}$ , the scalar product  $(f, g)$  is a real or complex number. Let  $\mathbb{C}$  denote the field of complex numbers. The map  $\{f, g\} \rightarrow (f, g)$  has the following properties:

(H1)  $(\alpha f, g) = \alpha(f, g)$   $(f, \alpha g) = \alpha^*(f, g)$ ,  $\alpha \in \mathbb{C}$ ,  $f, g \in \mathcal{H}$  where  $\alpha^*$  is the complex conjugate of  $\alpha$ ;

(H2)  $(f, g_1 + g_2) = (f, g_1) + (f, g_2)$ ,  $\forall f, g_1, g_2 \in \mathcal{H}$ ;

(H3)  $(f, g) = (g, f)^*$ ,  $\forall f, g \in \mathcal{H}$ ;

(H4)  $\|f\|_{\mathcal{H}}^2 = (f, f)$ ,  $\forall f \in \mathcal{H}$ .

**Definition 2.1.6.** Let  $\mathcal{Z}$  be any Banach space and  $\phi$  a real valued function on  $\mathcal{Z}$ . The function  $\phi$  is said to be [37]:

i) lower semi continuous at  $x \in \mathcal{Z}$  if, for every sequence  $\{x_n\}$  converging to  $x$ ,

$$\phi(x) \leq \liminf_{n \rightarrow \infty} \phi(x_n),$$

ii) upper semi continuous at  $x$  if

$$\phi(x) \geq \limsup_{n \rightarrow \infty} \phi(x_n),$$

iii) lower or upper semi continuous on a set  $\Gamma \subset \mathcal{Z}$ , if the corresponding statements hold for all  $x \in \Gamma$ .

Let  $X$  be a Banach space with the first and second duals denoted by  $X^*$  and  $X^{**}$  respectively. A sequence  $x_n^* \in X^*$  is said to converge weakly to  $x^*$ , denoted by

$$x_n^* \longrightarrow x^*,$$

if, for every  $x^{**} \in X^{**}$ ,

$$x^{**}(x_n^*) \longrightarrow x^{**}(x^*)$$

and it is said to converge in the weak star topology to  $x^*$ , denoted by

$$x_n^* \xrightarrow{w^*} x^*,$$

if, for every  $x \in X$ ,

$$x_n^* \longrightarrow x^*(x).$$

Since every element of  $X$  induces a continuous linear functional on  $X^*$  through the relation  $\hat{x}(x^*) \equiv x^*(x)$ , the canonical embedding  $X \hookrightarrow X^{**}$  exists. Hence the weak star topology is weaker than the weak topology [37].

For simplicity let  $(\Sigma, d_\Sigma)$  denote a complete separable metric space (a Polish space), and  $(\Sigma, \mathcal{B}(\Sigma))$  the corresponding measurable space, in which  $\mathcal{B}(\Sigma)$  is the  $\sigma$ -algebra generated by open sets in  $\Sigma$ . The material presented below regarding different spaces and their duals can be generalized to locally compact separable metric spaces  $(\Sigma, d_\Sigma)$ .

Let  $\mathcal{X}_0 \triangleq C_0(\Sigma)$  denote the Banach space of continuous functions on  $\Sigma$  that vanish at infinity,  $\mathcal{X}_1 \triangleq BC(\Sigma)$  the Banach space of bounded continuous functions on  $\Sigma$ , and  $\mathcal{X}_2 \triangleq BM(\Sigma)$  the Banach space of bounded measurable functions on  $\Sigma$ , all equipped with the sup-norm. Clearly,  $\mathcal{X}_0 \subset \mathcal{X}_1 \subset \mathcal{X}_2$ .

It is known that the dual space  $\mathcal{X}_0^*$  is isometrically isomorphic to  $\mathcal{M}_{rca}(\Sigma)$ , the Banach space of finite signed Borel measures on  $(\Sigma, \mathcal{B}(\Sigma))$  (also known as Radon measures), the dual space  $\mathcal{X}_1^*$  is isometrically isomorphic to  $\mathcal{M}_{rba}(\Sigma)$ , the Banach space of finitely additive finite regular signed measures on  $(\Sigma, \mathcal{B}(\Sigma))$ , and the dual space  $\mathcal{X}_2^*$  is isometrically isomorphic to  $\mathcal{M}_{ba}(\Sigma)$ , the Banach space of finitely additive signed measures on  $(\Sigma, \mathcal{B}(\Sigma))$ . Note that when  $\Sigma$  is compact then  $\mathcal{X}_1^*$  is isometrically isomorphic, the Banach space of countably additive signed measures on  $(\Sigma, \mathcal{B}(\Sigma))$ [38].

**Definition 2.1.7.** A Banach space  $X$  is said to be reflexive, if  $X^{**} = X$  [37].

For  $1 < p < \infty$ , the  $L_p$  spaces are reflexive Banach spaces. Indeed, for  $(1/p) + (1/q) = 1$ ,  $(L_p)^* = L_q$  and  $(L_q)^* = L_p$ . Hence  $(L_p)^{**} = L_p$  and so these spaces are reflexive.

It is well known that a closed bounded subset of a finite dimensional space is compact. Through this is false in infinite dimensional spaces, there is a similar result with respect to weak topologies. This is presented in the next theorem.

**Theorem 2.1.8.** A closed bounded subset of a reflexive Banach space is weakly compact [37].



**Theorem 2.1.9.** Suppose  $f : \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$  is lower semi continuous satisfying the following conditions:

$$-\infty < f(x) \neq +\infty, \forall x \in \mathbb{R}^d, \text{ and } \lim_{\|x\| \rightarrow \infty} f(x) = +\infty. \quad (2.1)$$

Then  $f$  attains its minimum [37].

In a finite dimensional space the closed unit ball is compact. An analogous result in infinite dimensional setting is Alaoglu's theorem.

**Theorem 2.1.10.** (Alaoglu's Theorem) The unit ball  $B_1(X^*)$  of the dual  $X^*$  of the Banach space  $X$  is weak star compact [37].

**Definition 2.1.11.** (Convex Set) A set  $S$  in a vector space  $X$  is called a convex set if the line segment joining any pair of points of  $S$  lies entirely in  $S$ . The former statement is equivalent to saying that for any pair of vectors  $u \in S, v \in S$ , the vector  $(1-t)u + tv \in S, \forall t \in [0, 1]$ .

**Definition 2.1.12.** (Convex Function) Let  $E$  be any Banach space and  $f$  a (possibly extended) real valued function defined on  $E$ . The function  $f$  is said to be convex if for every  $x, y \in E$  and  $\alpha \in [0, 1]$

$$f((1-\alpha)x + \alpha y) \leq (1-\alpha)f(x) + \alpha f(y).$$

Similarly, if  $E$  is replaced by a closed convex subset  $\Gamma$  on  $E$  and  $f$  satisfies the above inequality for all  $x, y \in \Gamma$ , then  $f$  is convex on  $\Gamma$ .

**Definition 2.1.13.** (Gateaux differential sub differential) A real valued functional  $f$  defined on  $E$  is said to be Gateaux differentiable at the point  $x_0 \in E$  if, for every  $y \in E$ , the limit

$$\lim_{s \rightarrow 0} (1/s) \{f(x_0 + sy) - f(x_0)\} \equiv df(x_0, y)$$

exists. Further, if  $y \rightarrow df(x_0, y)$  is continuous and linear, then there exists an  $e^* \in E^*$  in the dual space of  $E$ , dependent on  $x_0$ , such that

$$df(x_0, y) = (e^*, y).$$

The element  $e^* \in E^*$  satisfying the preceding identity, is called Gateaux gradient of  $f$  at  $x_0 \in E$ . The function  $f$  is said to be Gateux differentiable if it is so at every point  $x_0 \in E$  [37].

**Theorem 2.1.14.** A real valued function  $f$  defined on a Banach space  $E$  is weakly lower semi continuous if it is convex and continuously (linearly) Gateux differentiable.

## 2.1.2 Minimax Theory

Minimax techniques are often used in decision theory, game theory, statistics, based on the philosophy of minimizing the maximum possible loss. It started from two-player zero-sum game theory, covering both the cases where players take alternate moves and those where they make simultaneous moves. It has also been extended to general decision making in the presence of uncertainty.

Below the main minsup theorem is stated, which is used in this thesis; it is based on a generalization of the von Neumann's minimax theorem.

**Theorem 2.1.15.** *Let  $X$  be a compact Hausdorff space and  $Y$  an arbitrary set (not topologized). Let  $f$  be a real-valued function on  $X \times Y$  such that, for every  $y \in Y$ ,  $f(x, y)$  is lower semicontinuous on  $X$ . If  $f$  is convex on  $X$  and concave on  $Y$ , then there exists an  $x^* \in X$  such that*

$$\min_{x \in X} \sup_{y \in Y} f(x, y) = \sup_{y \in Y} f(x^*, y) = \sup_{y \in Y} \min_{x \in X} f(x, y). \quad (2.2)$$

*If in addition  $Y$  is a compact Hausdorff space for every  $x \in X$ ,  $f(x, y)$  is upper semicontinuous on  $Y$ , then there exists an  $(x^*, y^*) \in X \times Y$  saddle point, and*

$$\min_{x \in X} \max_{y \in Y} f(x, y) = f(x^*, y^*) = \max_{y \in Y} \min_{x \in X} f(x, y). \quad (2.3)$$

*Proof.* See [39]. □

The following theorem will be invoked to prove the equivalence between constrained and unconstrained optimization problems.

**Theorem 2.1.16.** *(Lagrange Duality) [40] (page 224-225). Let  $f$  be a real-valued convex functional defined on a convex subset  $\Omega$  of a vector space  $\mathcal{X}$  and let  $G$  be a convex mapping from  $\mathcal{X}$  into a normed space  $\mathcal{Z}$ . Suppose there exists an  $x_1 \in \mathcal{X}$  such that  $G(x_1) \prec 0$  (here  $\prec$  is used for the ordered vector space  $(\mathcal{Z}, \prec)$ ) and  $\inf\{f(x) : G(x) \preceq 0, x \in \Omega\}$  is finite, then*

$$\begin{aligned} & \inf\{f(x) : G(x) \preceq 0, x \in \Omega\} \\ &= \max_{z^* \succeq 0} \left( \inf_{x \in \Omega} \{f(x) + (G(x), z^*)\} \right) \end{aligned} \quad (2.4)$$

*and the maximum on the right is achieved by some  $z_0^* \succeq 0, z_0^* \in \mathcal{Z}$ .*

### 2.1.3 Measurable Space and Probability Space

The mathematical model for a random experiment is the probability space. Before it is introduced, the class of measurable functions which play a fundamental role in integration theory, and hence in measure theory are reviewed. There is an analogy between the concepts of topological space, open set, and continuous function and measurable space, measurable set, and measurable function.

**Definition 2.1.17.** (Topological Space) Let  $\mathcal{Z}$  be a set and  $\mathcal{B}_{\mathcal{Z}}$  a collection of subsets of  $\mathcal{Z}$ . Then  $\mathcal{B}_{\mathcal{Z}}$  is called a topology in  $\mathcal{Z}$  if the following properties hold [41].

- i)  $\emptyset \in \mathcal{B}_{\mathcal{Z}}$  and  $\mathcal{Z} \in \mathcal{B}_{\mathcal{Z}}$ ;
- ii) If  $Z_i \in \mathcal{B}_{\mathcal{Z}}, i = 1, 2, \dots, n$ , then  $\bigcap_{i=1}^n Z_i \in \mathcal{B}_{\mathcal{Z}}$
- iii) If  $\{Z_i\}$  is an arbitrary collection of elements of  $\mathcal{Z}$  (finite, countable, or uncountable), then  $\bigcup_i Z_i \in \mathcal{B}_{\mathcal{Z}}$ .

The pair  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}})$  is called a topological space and the members of  $\mathcal{B}_{\mathcal{Z}}$  are called open sets in  $\mathcal{Z}$ . If  $f : (\mathcal{Z}, \mathcal{B}_{\mathcal{Z}}) \rightarrow (Y, \mathcal{B}_Y)$ , then  $f$  is continuous provided  $f^{-1}(Y_i) \in \mathcal{B}_{\mathcal{Z}}$  is an open set for every open set  $Y_i \in \mathcal{B}_Y$ . Moreover,  $f$  is continuous at the point  $x_0 \in \mathcal{Z}$  if for every neighborhood (nbh)  $A$  of  $f(x_0)$  there exists a nbh  $B$  of  $x_0$  such that  $f(B) \subset A$ .

A topological space  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}})$  is said to be  $T_2$  (or is said to satisfy the  $T_2$  axiom) if given distinct  $x, y \in \mathcal{Z}$ , there exist disjoint open set  $U, V \in \mathcal{B}_{\mathcal{Z}}$  (that is,  $U \cap V = \emptyset$ ) such that  $x \in U$  and  $y \in V$ . A  $T_2$  space is also known as a Hausdorff space. A Hausdorff topology for a set  $\mathcal{Z}$  is a topology  $\mathcal{B}_{\mathcal{Z}}$  such that  $(\mathcal{Z}, \mathcal{B}_{\mathcal{Z}})$  is a Hausdorff space.

**Definition 2.1.18.** (Algebra) Let  $\Omega$  be a set of elementary outcomes and  $\mathcal{F}$  be a non-empty collection of subsets of  $\Omega$ . Then  $\mathcal{F}$  is called an Algebra on  $\Omega$  if the following properties hold [41].

- i)  $\Omega \in \mathcal{F}$  (The Sample Space is an element of  $\mathcal{F}$ );
- ii) If  $A \in \mathcal{F}$  then  $A^c = \Omega - A \in \mathcal{F}$ , where  $A^c$  is the complementation of  $A$  relative to  $\Omega$  (if a subset of  $\Omega$  belongs to  $\mathcal{F}$ , then so is its complement);
- iii) If  $A_i \in \mathcal{F}, i = 1, 2, \dots, n$ , then  $\bigcup_{i=1}^n A_i \in \mathcal{F}$  (if a finite number of subsets belong to  $\mathcal{F}$ , then so is their union).

Clearly, an algebra is a collection of subsets of a set  $\Omega$ , which a) contains  $\Omega$  and b) is closed under complementation and finite unions. The members of  $\mathcal{F}$  are called  $\mathcal{F}$ -measurable sets or measurable sets.

The collection of finite unions of half open intervals  $(a, b]$ ,  $-\infty < a < b \leq \infty$  in  $\mathbb{R}$  is considered. This collection is a field. However, the open interval  $(0, 1) = \bigcup_{n=1}^{\infty} (0, 1 - \frac{1}{n}]$  is not in the collection although it contains each interval  $(0, 1 - \frac{1}{n}]$ . Similarly, it does not contain the singletons  $\{x\}$ , although  $\{x\} = \bigcap_{n=1}^{\infty} (x - \frac{1}{n}, x]$ . Therefore, in order to consider sequences of events and convergence of sequence of events, it is necessary to extend the operations on events to countable set operations. This gives rise to a  $\sigma$ -algebra which is closed under countable unions.

**Definition 2.1.19.** ( $\sigma$ -Algebra) An algebra  $\mathcal{F}$  on  $\Omega$  is called a  $\sigma$ -Algebra on  $\Omega$  if it is closed under countable unions, that is if the following properties hold [41].

- i)  $\Omega \in \mathcal{F}$  ;
- ii) If  $A \in \mathcal{F}$  then  $A^c \in \mathcal{F}$ ;
- iii) If  $A_i \in \mathcal{F}$ ,  $i = 1, 2, \dots$  then  $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$ .

If  $\mathcal{F}$  is a field the pair  $(\Omega, \mathcal{F})$  is called a measurable space and the elements of  $\mathcal{F}$  are called Events and are said to be measurable sets in  $\Omega$ . Fields and  $\sigma$ -fields are convenient mathematical objects which express how much is known about the outcome of the experiment.

**Remark 2.1.20.** If  $f : X \rightarrow Y$ ,  $X$  is measurable Space,  $Y$  is topological Space, then  $f$  is said to be measurable if  $f^{-1}(V)$  is a measurable set in  $X$  for every open set  $V$  in  $Y$ .

Since the intersection of arbitrary  $\sigma$ -algebras of subset of  $\Omega$  is a  $\sigma$ -algebra of subsets of  $\Omega$ , then for an arbitrary family  $\mathcal{A}$  of subsets of  $\Omega$  there is a smallest  $\sigma$ -algebra  $\mathcal{F}$  in  $\Omega$  such that  $\mathcal{A} \subset \mathcal{F}$ .

**Theorem 2.1.21.** (Smallest  $\sigma$ -algebra) Let  $\Omega$  be a sample space and  $\mathcal{A}$  be a collection of subsets of  $\Omega$ . There exists a smallest  $\sigma$ -algebra  $\mathcal{F}(\mathcal{A})$  on  $\Omega$  containing  $\mathcal{A}$ , which is constructed by

$$\mathcal{F}(\mathcal{A}) = \bigcap_i \{ \mathcal{N}_i; \mathcal{N}_i \text{ is a } \sigma\text{-algebra on } \Omega, \mathcal{A} \subset \mathcal{N}_i \}.$$

This is called the  $\sigma$ -algebra generated by  $\mathcal{A}$ , and it is often denoted by  $\mathcal{F}(\mathcal{A}) = \sigma(\mathcal{A})$ .

**Borel Set.** Let  $X$  be a topological space (e.g.,  $X = \mathfrak{R}^n$ ; the collection of all  $n$ -tuples  $\{x = (x_1, x_2, \dots, x_n) : x_i \in \mathfrak{R}, 1 \leq i \leq n\}$ ). Then there exists a smallest  $\sigma$ -algebra  $\mathcal{F}$  on  $X$  such that every open set  $\mathcal{A} \subset X$  belongs to  $\mathcal{F}$ . The elements  $A \in \mathcal{F}$  called Borel sets and the  $\sigma$ -algebra  $\mathcal{F} = \mathcal{F}(X)$  is called a Borel  $\sigma$ -algebra. For example, if  $X = \mathfrak{R}^n$ , and  $\mathcal{A}$  is the collection of all open sets of  $\mathfrak{R}^n$ , the Borel  $\sigma$ -algebra denoted by  $\mathcal{B}(\mathfrak{R}^n)$ , contains all open sets, their complements (closed sets), all the countable unions of open sets, and all the countable unions of closed sets. In fact,  $\mathcal{B}(\mathfrak{R}^n)$  = the smallest  $\sigma$ -algebra of subsets of  $\mathfrak{R}^n$  containing all sets of the form  $\{x : x_1 \in A_1, x_2 \in A_2, \dots, x_n \in A_n\}$ , where  $A_j$  are intervals on  $\mathfrak{R}$  which are closed, open, semi-open, points, etc. Clearly,  $\mathcal{A}$  = Collection of all open intervals of  $\mathfrak{R}^n$  is not a  $\sigma$ -algebra, but there exists many  $\sigma$ -algebra containing  $\mathcal{A}$  as a subset. The smallest  $\sigma$ -algebra containing  $\mathcal{A}$  is the  $\sigma$ -algebra generated by  $\mathcal{A}$ . The pair  $(\mathfrak{R}^n, \mathcal{B}(\mathfrak{R}^n))$  is a measurable space, called, the Borel measurable space.

**Probability Space.** In order to grade the possibility of occurrences of events associated with a random experiment a function (a map) has to be defined which attaches a numerical value to events  $A \in \mathcal{F}$ . A function  $\mu : \mathcal{F} \rightarrow [0, \infty]$  is called a finite-additive set function, if  $\mu$  satisfies the following two conditions.

$$(FA1) \quad \mu(\phi) = 0;$$

$$(FA2) \quad \mu(A \cup B) = \mu(A) + \mu(B), \text{ if } A, B \in \mathcal{F} \text{ and } A \cap B = \emptyset.$$

A finite-additive set function  $\mu$  on an algebra  $\mathcal{F}$  (or a  $\sigma$ -algebra) is called a measure, if it is countably-additive and a probability measure if it is countably-additive and  $\mu(\Omega) = 1$ , hence the following definition.

**Definition 2.1.22.** (Probability Measure) Let  $(\Omega, \mathcal{F})$  be a measurable space. The map

$$P : \mathcal{F} \rightarrow [0, 1], \quad P(A) \in [0, 1], \quad \forall A \in \mathcal{F}$$

is called a probability measure on  $(\Omega, \mathcal{F})$  if it satisfies the following properties.

- i)  $P(\phi) = 0$ ;
- ii)  $P(\Omega) = 1$ ;
- iii)  $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$ , if  $A_i \in \mathcal{F}, \forall i$  and  $\{A_j\}_{j=1}^{\infty}$  are disjoint, e.g.,  $A_i \cap A_j = \emptyset, \forall i \neq j$ .

The triple  $(\Omega, \mathcal{F}, P)$  is called a probability space [41].

**Completeness of Probability Space.** A probability Space  $(\Omega, \mathcal{F}, P)$  is said to be Complete if whenever  $B \in \mathcal{F}$  and  $P(B) = 0$  then  $A \in \mathcal{F}$  for all  $A \subset B$ . The subsets  $A$  of an event  $B$  of zero probability is called a null set, therefore  $(\Omega, \mathcal{F}, P)$  is complete if  $\mathcal{F}$  includes all events of zero probability. Any probability space  $(\Omega, \mathcal{F}, P)$  which is not complete can be uniquely extended to the  $\sigma$ -algebra  $\overline{\mathcal{F}} = \mathcal{F} \vee \{\text{Null Sets}\}$ .

**Definition 2.1.23.** Let  $(\Omega, \mathcal{F})$  denote a measurable space and  $P$  a positive measure  $\mu$  on  $\Omega$ . Let  $f$  be a measurable function on  $(\Omega, \mathcal{F})$ . Define

$$\|f\|_p \triangleq \left\{ \int_{\Omega} |f|^p(\omega) d\mu(\omega) \right\}^p, \quad 1 \leq p < \infty.$$

$L^p(\Omega, \mathcal{F}, P)$  is the set of all measurable functions  $f$  on  $(\Omega, \mathcal{F})$  for which  $\|f\|_p < \infty$ .  $\|f\|_p$  denotes the  $L^p$ -norm of  $f$ .

**Definition 2.1.24.** (Mutually Independence) The events  $\{A_i\}_{i=1}^n$  are said to be mutually independent if

$$P(A_{i_1} \cap A_{i_2} \cdots \cap A_{i_k}) = P(A_{i_1})P(A_{i_2}) \cdots P(A_{i_k})$$

for all non-empty subsets  $\{i_1, i_2, \dots, i_k\}$  of  $\{1, 2, \dots, n\}$  [41].

**Remark 2.1.25.** Pairwise independence between events does not imply mutual independence of all events. Disjoint events are not independent since if  $A, B$  are disjoint and independent events then  $0 = P(A \cap B) = P(A)P(B) \Rightarrow P(A)$  and/or  $P(B)$  are zero.

## 2.1.4 Random Variables

Let  $(\Omega_1, \mathcal{F}_1)$  and  $(\Omega_2, \mathcal{F}_2)$  be two measurable spaces, and let  $f : (\Omega_1, \mathcal{F}_1) \rightarrow (\Omega_2, \mathcal{F}_2)$ . Then the function  $f$  is called  $\mathcal{F}_1/\mathcal{F}_2$  or  $\mathcal{F}_1$  measurable if

$$f^{-1}(A) \triangleq \{\omega : f(\omega) \in A\} \in \mathcal{F}_1, \quad \forall A \in \mathcal{F}_2.$$

The set  $f^{-1}(A)$  is called the inverse image of  $A \in \mathcal{F}_2$ . If  $f : \Omega \rightarrow Y$  where  $(\Omega, \mathcal{F})$  is a measurable space,  $Y$  is a topological space (e.g.,  $\mathbb{R}^n$ ), then  $f$  is  $\mathcal{F}/\mathcal{B}(\mathbb{R}^n)$ -measurable provided  $f^{-1}(V) \in \mathcal{F}$  for every open set  $V \subset Y$ .

The  $\sigma$ -algebra  $\mathcal{F}(f)$  generated by  $f$  is the smallest  $\sigma$ -algebra on  $\Omega$  containing all the sets  $\{f^{-1}(V) : V \subset Y \text{ is open}\}$  and  $f$  will be  $\mathcal{F}(f)/\mathcal{B}(\mathbb{R}^n)$ . Moreover, if  $Y = \mathbb{R}^n$  then

$$\mathcal{F}(f) = \{f^{-1}(V) : V \in \mathcal{B}(\mathbb{R}^n)\}.$$

Clearly, if  $(\Omega, \mathcal{B})$  is a Borel measurable space and  $f : \Omega \rightarrow Y$ , where  $Y$  is a topological space and  $f$  is a continuous function, then from the definition of continuous function

$$f^{-1}(V) \in \mathcal{B}, \quad \forall \text{ open set } V \subset Y. \quad (2.5)$$

Hence, every continuous function is Borel measurable, called Borel function, e.g.,  $f : (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n)) \rightarrow (\mathbb{R}^m, \mathcal{B}(\mathbb{R}^m))$  is a Borel function.

If a probability measure  $P$  on  $(\Omega, \mathcal{F})$  is defined, where  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$ , is a measurable function then  $X$  is called a Random Variable (RV) defined on the probability space  $(\Omega, \mathcal{F}, P)$ .

**Definition 2.1.26.** (Random Variable) Let  $X : \Omega \rightarrow \mathbb{R}^n$  be a function defined on a probability space  $(\Omega, \mathcal{F}, P)$ . Then  $X$  is called an  $n$ -dimensional Random Variable (measurable function)

$$X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$$

if for every  $A \in \mathcal{B}(\mathbb{R}^n)$  the set

$$X^{-1}(A) \triangleq \{\omega : X(\omega) \in A\} \in \mathcal{F}.$$

Clearly, the  $\sigma$ -algebra  $\mathcal{F}^X$  (or  $\mathcal{F}(X)$ ) generated by  $X$  is the smallest  $\sigma$ -algebra on  $\Omega$  containing all the sets

$$X^{-1}(A) : A \subset \mathbb{R}^n \text{ is open}$$

under which  $X$  is measurable. Equivalently,

$$\mathcal{F}^X = X^{-1}(\mathcal{B}(\mathbb{R}^n)) = \{X^{-1}(B) : B \in \mathcal{B}(\mathbb{R}^n)\}$$

is the smallest  $\sigma$ -algebra on  $\Omega$  under which  $X$  is measurable. However,  $\mathcal{B}(\mathbb{R}^n)$  is generated by products of open sets of the form

$$\{x : -\infty < x_1 \leq \alpha_1, \dots, -\infty < x_n \leq \alpha_n\}, \quad \alpha_j \in \mathbb{R}, 1 \leq j \leq n \quad (2.6)$$

therefore for  $X$  to be a RV it is sufficient for every set of the form

$$\{\omega : X_1(\omega) \leq \alpha_1, \dots, X_n(\omega) \leq \alpha_n\}, \alpha_j \in \mathfrak{R}, 1 \leq j \leq n$$

to be an event. This is because  $\mathcal{B}(\mathfrak{R}^n)$  is the family of subsets obtained by starting with (2.6) and taking repeatedly all complements, countable unions, intersections. Also, if  $\{X_t : 0 \leq t \leq T\}$  is a family of random variables  $\mathcal{F}_t^X \triangleq \sigma(X_t : 0 \leq t \leq T) = \bigvee_{t \in T} \mathcal{F}(X_t) = \sigma(\bigcup_{t \in T} \mathcal{F}(X_t))$  is the smallest  $\sigma$ -algebra on  $\Omega$  under which  $\{X_t : 0 \leq t \leq T\}$  are measurable.

**Complex Random Variables.** A complex square matrix  $Q$  is called Hermitian if  $Q = Q^\dagger$  (where  $\dagger$  denotes complex conjugate transpose), and has the following properties [42]:

- i) The eigenvalues of a Hermitian matrix are real.
- ii) The diagonal elements of a Hermitian matrix are real.
- iii) The complex conjugate of a Hermitian matrix is a Hermitian matrix.
- iv) If  $Q$  is a Hermitian matrix, and  $B$  is a complex matrix of same order as  $Q$ , then  $BQB$  is a Hermitian matrix.
- v) A matrix is symmetric if and only if it is real and Hermitian.
- vi) Hermitian matrices are a vector subspace of the vector space of complex matrices. The real symmetric matrices are a subspace of the Hermitian matrices.
- vii) Hermitian matrices are also called self-adjoint since if  $Q$  is Hermitian, then in the usual inner product of  $\mathbb{C}^n$ , we have  $(uQv) = (Quv)$  for all  $u, v \in \mathbb{C}^n$ .

A complex Random Variable (RV)  $Z \in \mathbb{C}^n$  is simply a pair of real RVs of  $\mathfrak{R}^n$  such that

$$Z = X + jY. \tag{2.7}$$

It is therefore always possible to treat all the problems concerning complex RVs by using a real RV of  $\mathfrak{R}^{2n}$  dimension.

A complex random vector  $Z \in \mathbb{C}^n$  is said to be Gaussian if the real random vector  $\bar{Z} \in \mathfrak{R}^{2n}$  consisting of its real and imaginary parts,  $\bar{Z} = \begin{bmatrix} \Re(Z) \\ \Im(Z) \end{bmatrix}$  is Gaussian. Thus,



to specify the distribution of a complex Gaussian random vector  $Z$ , it is necessary to specify the expectation and covariance of  $\bar{Z}$ , namely

$$E[\bar{Z}] \in \mathfrak{R}^{2n} \quad \text{and} \quad E[(\bar{Z} - E[\bar{Z}])(\bar{Z} - E[\bar{Z}])^T] \in \mathfrak{R}^{2n \times 2n}.$$

A complex Gaussian random vector  $Z$  is circularly symmetric if the covariance of the corresponding  $\bar{Z}$  has the structure

$$E[(\bar{Z} - E[\bar{Z}])(\bar{Z} - E[\bar{Z}])^T] = \begin{bmatrix} \Re(Q) & -\Im(Q) \\ \Im(Q) & \Re(Q) \end{bmatrix} \quad (2.8)$$

for some Hermitian non-negative definite  $Q \in \mathbb{C}^{n \times n}$ . Note that the real part of an Hermitian matrix is symmetric and the imaginary part of an Hermitian matrix is anti-symmetric and thus the matrix appearing in (2.8) is real and symmetric. In this case  $E[(\bar{Z} - E[\bar{Z}])(\bar{Z} - E[\bar{Z}])^T] = Q$ , and thus, a circularly symmetric complex Gaussian random vector  $Z$  is specified by prescribing  $E[Z]$  and  $E[(Z - E[Z])(Z - E[Z])^\dagger]$ .

### 2.1.5 Distribution Function

Let  $(\Omega, \mathcal{F}, P)$  be a probability space and let  $X : (\Omega, \mathcal{F}) \rightarrow (\Omega_1, \mathcal{F}_1)$  be an  $\mathcal{F}/\mathcal{F}_1$ -measurable RV. From the point of view of computations, it is often convenient to work with an induced measure on  $\mathcal{F}_1$ . This amounts to defining the probability measure induced by the RV on its range rather than treat points with respect to the measure  $P$ , and work with a probability measure on  $\mathcal{F}_1$  with  $\omega \in \Omega_1$  as its sample values.

For the specific case, the RV,  $X : (\Omega, \mathcal{F}) \rightarrow (\Omega_1, \mathcal{F}_1)$  induces a probability measure  $P_X$  on  $(\Omega_1, \mathcal{F}_1)$  by

$$\begin{aligned} P_X(A_1) &\triangleq P \circ X^{-1}(A_1) = P(\{\omega : X(\omega) \in A_1\}) \\ &= P(X \in A_1), \quad A_1 \in \mathcal{F}_1. \end{aligned}$$

If  $(\Omega_1, \mathcal{F}_1) = (\mathfrak{R}, \mathcal{B}(\mathfrak{R}))$  then someone can work with a probability measure on  $\mathcal{B}(\mathfrak{R})$  with  $x \in \mathfrak{R}$  as its sample points.

**Definition 2.1.27.** (*Probability Distribution*) Let  $(\Omega, \mathcal{F}, P)$  be a Probability Space and  $X : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}, \mathcal{B}(\mathfrak{R}))$  a RV. The function  $F_X(\cdot)$  defined as

$$F_X(x) \triangleq P(\{\omega : X(\omega) \leq x\}) = P_X(X \leq x)$$

is called the (cumulative) probability distribution of  $X$ .

Thus, the relationship

$$P_X(A) = P(\{\omega : X(\omega) \in A\})$$

defines a probability measure  $P_X$  on  $(\mathfrak{R}, \mathcal{B}(\mathfrak{R}))$ . Note that  $F_X(x)$  is a probability distribution defined on  $\mathfrak{R}$ , e.g., it corresponds to the probability measure corresponding to  $P$  induced by  $X(\cdot)$  on  $\mathfrak{R}$ .

Suppose  $X_1, X_2, \dots, X_n$  are  $n$  real-valued RV's and  $X = (X_1, X_2, \dots, X_n)$ , then

$$X : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}^n, \mathcal{B}(\mathfrak{R}^n))$$

is a measurable function. The function

$$F_X(x) = F_X(x_1, x_2, \dots, x_n) = P(\{\omega : X_i(\omega) \leq x_i, i = 1, \dots, n\}), x \in \mathfrak{R}^n$$

is called the joint probability distribution function of  $X$ . Similarly as above, the relationship

$$P_X(A) = P(\{\omega : X(\omega) \in A\}), \quad A \in \mathcal{B}(\mathfrak{R}^n)$$

defines a Borel probability measure.

A real-valued RV  $X$  is said to be discrete if there exists a countable set  $S = \{x_i\}$  such that

$$\sum_{x_i \in S} P(\{\omega : X(\omega) = x_i\}) = 1.$$

If  $X$  is discrete, then the distribution function  $F_X$  is a function which is constant except for jumps at  $x_i, i = 1, 2, \dots$ , the size of the jump at  $x_i$  being  $P(\{\omega : X(\omega) = x_i\})$ . For an arbitrary Borel set  $A$ ,

$$P_X(A) = \sum_{x_i \in A \cap S} P(\{\omega : X(\omega) = x_i\}).$$

Let  $P$  be a probability measure on  $(\mathfrak{R}^n, \mathcal{B}(\mathfrak{R}^n))$ . It is said to be singular (with respect to the Lebesgue measure) if there exist a set  $S \in \mathcal{B}(\mathfrak{R}^n)$  such that  $P(S) = 1$  and the Lebesgue measure of  $S$  is zero. On the other hand,  $P$  is said to be absolutely

continuous (w.r.t. the Lebesgue measure) if Lebesgue measure of  $(A)$  implies  $P(A) = 0$ . Clearly, if  $X_1, X_2, \dots, X_n$  are discrete RV's, then  $P_X$  is singular. If  $X_1, X_2, \dots, X_n$  are such that  $P_X$  is absolutely continuous, then there exists a non-negative Borel function  $p_X(x), x \in \mathfrak{R}^n$  such that

$$P_X(A) = \int_A p_X(x) dx, A \in \mathcal{B}(\mathfrak{R}^n).$$

The function  $p_X$  is called the probability density function for  $\mathbf{X}$ . In terms of the distribution

$$F_X(x_1, x_2, \dots, x_n) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} p_X(x_1, \dots, x_n) dx_1 \dots dx_n$$

which implies

$$p_X(x_1, \dots, x_n) \triangleq \frac{\partial^n}{\partial x_1 \dots \partial x_n} F_X(x_1, x_2, \dots, x_n).$$

**Definition 2.1.28.** (Stochastic Kernel) Given a measurable space  $(\Omega, \mathcal{F})$  on which the RVs,  $X$  and  $Y$  are defined, via  $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \Sigma_X)$ , and  $Y : (\Omega, \mathcal{F}) \rightarrow (\mathcal{Y}, \Sigma_Y)$ , respectively, then the relation between the RV  $X$  and the RV  $Y$  is defined via a probabilistic mapping. The mapping  $\mu : \mathcal{X} \times \Sigma_Y \rightarrow [0, 1]$  satisfies the following two conditions:

- (i) For every  $x \in \mathcal{X}$ , the set function  $\mu(x, \cdot)$  is a probability measure on  $\Sigma_Y$  (possibly finite additive);
- (ii) For every  $F \in \Sigma_Y$ , the function  $\mu(\cdot, F)$  is  $\mathcal{X}$ -measurable.

The mapping  $\mu$  is called a stochastic kernel or transition probability.

### 2.1.6 Duality Relation Between KL Distance and Free Energy

Next the basic definitions and duality relations are introduced between Kullback-Leibler (KL) distance, free energy, and cumulant moment generating function.

**Definition 2.1.29.** Let  $\nu, \mu \in \mathcal{M}_1(\Sigma)$  (the set of probability measures) and  $\ell : \Sigma \rightarrow \mathfrak{R}$  a measurable function.

- 1) The moment generating function of  $\ell$  with respect to  $\mu$  is defined by

$$M_\mu(s) \triangleq E_\mu(e^{s\ell}) = \int_\Sigma e^{s\ell} d\mu \in (0, \infty], s \in \mathfrak{R}. \quad (2.9)$$

2) The cumulant generating function of  $\ell$  with respect to  $\mu$  is defined by

$$\Psi_\mu(s) \triangleq \log M_\mu(s) = \log \int_\Sigma e^{s\ell} d\mu, \quad s \in \mathfrak{R}. \quad (2.10)$$

3) The free energy of  $\ell$  with respect to  $\mu$  is defined by  $\mathcal{E}(\ell, \mu) \triangleq \Psi_\mu(1) \in (-\infty, \infty]$ .

4) The KL distance of  $\nu \in \mathcal{M}_1(\Sigma)$  with respect to  $\mu \in \mathcal{M}_1(\Sigma)$  is defined by

$$H(\nu|\mu) \triangleq \begin{cases} \int_\Sigma \log\left(\frac{d\nu}{d\mu}\right) d\nu, & \text{if } \nu \ll \mu \\ +\infty, & \text{otherwise.} \end{cases}$$

It can be shown that  $\mathcal{E}(\ell, \mu)$  as a function of  $\ell$  is convex,  $H(\nu|\mu)$  as a function of  $\mu, \nu \in \mathcal{M}(\Sigma)$  is convex in both arguments,  $M_\mu(s), \Psi_\mu(s)$  are convex functions of  $s \in \mathfrak{R}$ , and  $H(\nu|\mu) \geq 0$ , and  $H(\nu|\mu) = 0$ , if and only if  $\mu = \nu$ . Moreover,  $H(\nu|\mu)$  is often used as a measure of discrepancy between two probability measures.

The moment generating function (2.9) and cumulant general function (2.10) are often employed as pay-off functions in nonlinear stochastic control problems to achieve robustness. Such pay-off's are called risk-sensitive [43],[44].

## 2.2 Change of Probability Measure

A basic technique used throughout this dissertation is a change of probability measure starting with probability  $P$ . A new probability  $\bar{P}$  is defined such that under  $\bar{P}$  the observations are independent and identically distributed random variables. Calculations take place in the mathematically ideal world of  $\bar{P}$  which allows interchange of expectations and summations. They are then related to the real world by an inverse change of measure. The measure change concept is the key to many of the results in the following Chapters.

Change of measure is a fundamental theorem in measure theory known as Radon-Nikodym Theorem. A version of Radon-Nikodym Theorem suitable for probability measures, is stated here.

**Theorem 2.2.1.** (Radon-Nikodym Theorem)

Let  $(\Omega, \mathcal{F}, P)$  be a probability space, and let  $\bar{P}$  be another measure defined also on  $\mathcal{F}$  such that  $\bar{P}$  is absolutely continuous with respect to  $P$  ( $\bar{P} \ll P$ ), namely,

$$P(B) = 0 \Rightarrow \bar{P}(B) = 0, \quad \forall B \in \mathcal{F}. \quad (2.11)$$

Then there exists an  $\mathcal{F}$ -measurable function  $\phi : \Omega \rightarrow \mathfrak{R}$ , such that  $\phi \in L^1(\Omega, \mathcal{F}, P)$  and,

$$\bar{P}(B) = \int_B \phi(\omega) dP(\omega), \quad \forall B \in \mathcal{F}. \quad (2.12)$$

The function  $\phi$  is unique except on a subset of  $P$ -measure zero.

This function  $\phi$  is often written as  $\phi = \left. \frac{d\bar{P}}{dP} \right|_{\mathcal{F}}$  and is called the Radon-Nikodym derivative (RND) since it satisfies

$$\bar{P}(B) = \int_B d\bar{P} = \int_B \phi dP, \quad \forall B \in \mathcal{F}. \quad (2.13)$$

**Lemma 2.2.2.** Let  $(\Omega, \mathcal{F})$  be a measurable space, and let  $P$ , and  $\bar{P}$  two measures defined on  $(\Omega, \mathcal{F})$  such that  $P$  is absolutely continuous with respect to  $\bar{P}$  and vice versa (mutually absolutely continues), namely,

$$\begin{aligned} P(B) = 0 &\Rightarrow \bar{P}(B) = 0, \quad \forall B \in \mathcal{F}, \\ \bar{P}(B) = 0 &\Rightarrow P(B) = 0, \quad \forall B \in \mathcal{F}. \end{aligned} \quad (2.14)$$

Then, given  $\phi(\omega) \triangleq \left. \frac{d\bar{P}}{dP} \right|_{\mathcal{F}}$  the two measures can be expressed as

$$\begin{aligned} \bar{P}(B) &= \int_B \phi(\omega) dP, \\ P(B) &= \int_B \phi^{-1}(\omega) d\bar{P}. \end{aligned} \quad (2.15)$$

Equivalently, for any random variable  $X : \Omega \rightarrow \mathfrak{R}$ , the following holds

$$\begin{aligned} \bar{E}[X] &= E[\phi X] = E\left[\frac{d\bar{P}}{dP} X\right], \\ E[X] &= \bar{E}[\phi^{-1} X] = \bar{E}\left[\frac{dP}{d\bar{P}} X\right] \end{aligned} \quad (2.16)$$

where  $\bar{E}$  and  $E$  denote expectations under  $\bar{P}$  and  $P$ , respectively.

$$\begin{aligned}\bar{P}(\omega) &= \int_{\omega} \phi(\omega) dP(\omega) = 1, \\ P(\omega) &= \int_{\omega} \phi^{-1}(\omega) d\bar{P}(\omega) = 1.\end{aligned}\tag{2.17}$$

**Theorem 2.2.3.** (Conditional Bayes Theorem)

Suppose  $(\Omega, \mathcal{F})$  be a measurable space and  $\mathcal{G} \subset \mathcal{F}$  is a sub- $\sigma$ -field. Suppose  $P$ , and  $\bar{P}$  are two measures defined on  $(\Omega, \mathcal{F})$  such that  $P$  is absolutely continuous with respect to  $\bar{P}$  and vice versa, with RND  $\frac{d\bar{P}}{dP} = \Lambda$  and  $\frac{dP}{d\bar{P}} = \Lambda^{-1}$ . Then if  $X$  is any integrable  $\mathcal{F}$ -measurable random variable,

$$\begin{aligned}\bar{E}[X|\mathcal{G}] &= \frac{E[\Lambda X|\mathcal{G}]}{E[\Lambda|\mathcal{G}]}, \quad \bar{P} - a.s. \\ E[X|\mathcal{G}] &= \frac{\bar{E}[\Lambda^{-1}X|\mathcal{G}]}{\bar{E}[\Lambda^{-1}|\mathcal{G}]}, \quad P - a.s.\end{aligned}\tag{2.18}$$

*Proof.* See [41]. □

## 2.2.1 Change of Probability Measure for Random Processes

**Definition 2.2.4.** A Random Process (RP)  $\{\Phi_k\}$ ,  $k \in N_0 \triangleq \{0, 1, 2, 3, \dots\}$  is said to be  $\{\mathcal{F}_k\}_{k \geq 0}$  adapted if  $\Phi_k$  is  $\mathcal{F}_k$ -measurable for every  $k \in N_0$ .

**Lemma 2.2.5.** Let  $(\Omega, \mathcal{F}, P)$  be a complete probability space and let  $\{\mathcal{F}_k\}$ ,  $\{\mathcal{F}_k^y\}$ ,  $k \in N_0$  be complete sub-sigma fields of  $\mathcal{F}$  such that  $\mathcal{F}_k^y \subset \mathcal{F}_k \subset \mathcal{F}$ ,  $\forall k \in N_0$ .

Let  $\Phi : \Omega \times N_0 \rightarrow \mathfrak{R}$ , a RP such that  $\{\Phi_k\}$  is  $\{\mathcal{F}_k\}$  adapted, and  $\Phi \in L^1(\Omega, \mathcal{F}, P)$ .

Then

$$E[\Phi_k|\mathcal{F}_k^y] = \frac{\bar{E}[\Lambda_k \Phi_k|\mathcal{F}_k^y]}{\bar{E}[\Lambda_k|\mathcal{F}_k^y]},\tag{2.19}$$

where  $\Lambda_k \triangleq \frac{dP}{d\bar{P}} \Big|_{\mathcal{F}_k}$ ,  $\Lambda_k^{-1} \triangleq \frac{d\bar{P}}{dP} \Big|_{\mathcal{F}_k}$   $P - a.s.$ ,  $E[\Lambda_k^{-1}] = 1$

### 2.2.2 Change of Probability Measure for Linear Systems

Let  $(\Omega, \mathcal{F}, P)$  be a complete probability space on which  $\{x_k\}, \{y_k\}, k \in N_0$  are defined by

$$x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1}, \quad x_0 \in \mathfrak{R}^n \quad (2.20)$$

$$y_k = C_k x_k + D_k v_k, \quad y_0 \in \mathfrak{R}^d. \quad (2.21)$$

Assume the following conditions hold:

- (i)  $w : \Omega \times N_0 \rightarrow \mathfrak{R}^m$  is an independent and identically distributed (iid) sequence with density

$$\Phi_{w_k}(w) = \frac{1}{(2\pi)^{m/2}} e^{-\frac{w^T w}{2}}, \quad w_k \sim N(0, I_m);$$

- (ii)  $v : \Omega \times N_0 \rightarrow \mathfrak{R}^d$  is an iid sequence with density

$$\Phi_{v_k}(v) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{v^T v}{2}}, \quad v_k \sim N(0, I_d);$$

- (iii)  $x_0 : \Omega \rightarrow \mathfrak{R}^n$  is a Random Variable (RV);

- (iv)  $\{x_0, \{w_k\}, \{v_k\}\}$  are mutually independent;

- (v)  $D_k > 0, \forall k \in N_0$ .

Next, the  $\sigma$ -algebras are defined

$$\begin{aligned} \mathcal{G}_k^0 &\triangleq \sigma\{x_0, x_1, \dots, x_k, y_0, y_1, \dots, y_k\}, \quad \forall k \in N_0 \\ \mathcal{F}_k^{0, \mathcal{Y}} &\triangleq \sigma\{y_0, y_1, \dots, y_k\}, \quad \forall k \in N_0. \end{aligned}$$

$\{\mathcal{G}_k\}, \{\mathcal{F}_k^{\mathcal{Y}}\}, \forall k \in N_0$  are the complete filtrations generated by  $\{x_0, x_1, \dots, x_k, y_0, y_1, \dots, y_k\}, \{y_0, y_1, \dots, y_k\}$ , respectively,  $\forall k \in N_0$ .

#### A. Measure Change of the Observation Process

The RV on  $(\{\mathcal{G}_k\}, P)$ , is considered, defined by

$$\lambda_k^{-1} \triangleq |D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)}, \quad k \in N_0 \quad (2.22)$$

setting also

$$\Lambda_k^{-1} \triangleq \prod_{s=0}^k \lambda_s^{-1}. \quad (2.23)$$

Then  $\Lambda_k^{-1} > 0$ , a.s. and in particular  $E[\Lambda_k^{-1}] = 1$ ,  $\forall k \in N_0$ , which can be shown as follows

$$\begin{aligned} E[\Lambda_k^{-1}] &= E\left[\prod_{s=0}^k |D_s| \frac{\Phi_{v_s}(y_s)}{\Phi_{v_s}(v_s)}\right] = E\left[E\left[\prod_{s=0}^k |D_s| \frac{\Phi_{v_s}(y_s)}{\Phi_{v_s}(v_s)} \middle| \mathcal{G}_{k-1}\right]\right] \\ &= E\left[\prod_{s=0}^{k-1} |D_s| \frac{\Phi_{v_s}(y_s)}{\Phi_{v_s}(v_s)} E\left[|D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} \middle| \mathcal{G}_{k-1}\right]\right]. \end{aligned}$$

But given

$$\prod_{s=0}^{k-1} |D_s| \frac{\Phi_{v_s}(y_s)}{\Phi_{v_s}(v_s)} = \Lambda_{k-1}^{-1}$$

and

$$E\left[|D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} \middle| \mathcal{G}_{k-1}\right] = \int_{\mathbb{R}^d} |D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} \Phi_{v_k}(v_k) dv_k = \int_{\mathbb{R}^d} |D_k| \Phi_{v_k}(y_k) |D_k|^{-1} dy_k = 1$$

(here change of variables is being used,  $v_k = D_k^{-1}(y_k - C_k x_k)$ ,  $dv_k = |D_k|^{-1} dy_k$ ).

Therefore

$$E[\Lambda_k^{-1}] = E[\Lambda_{k-1}^{-1}] = \dots = E[\Lambda_0^{-1}] = 1.$$

Hence,  $\Lambda_k^{-1}$  can be used to define a new probability measure  $\bar{P}$  on  $(\Omega, \{\mathcal{G}_k\})$  through the RND

$$\frac{d\bar{P}}{dP}\bigg|_{\mathcal{G}_k} = \Lambda_k^{-1} \text{ or } d\bar{P}(A) = \int_A \Lambda_k^{-1}(\omega) dP(\omega), \quad \forall A \in \mathcal{G}_k. \quad (2.24)$$

Next it is shown that under the new measure  $\bar{P}$ , the observation process  $\{y_k\}$ ,  $\forall k \in N_0$  is a sequence of iid RVs with density  $\{\Phi_{v_k}(\cdot)\}$ ,  $\forall k \in N_0$ .

**Lemma 2.2.6.** *Let  $(\Omega, \mathcal{F}, P)$  be a probability space on which (2.20), (2.21) are defined.*

*Next, define*

$$\frac{d\bar{P}}{dP}\bigg|_{\mathcal{G}_k} = \Lambda_k^{-1} = \prod_{s=0}^k |D_s| \frac{\Phi_{v_s}(y_s)}{\Phi_{v_s}(v_s)},$$

*where  $v_k = D_k^{-1}(y_k - C_k x_k)$ . Then under measure  $\bar{P}$  the sequence  $\{y_k\}$ ,  $k \in N_0$  is iid with density*

$$\Phi_{y_k}(y) = \frac{e^{-y^T y/2}}{(2\pi)^{d/2}} = \Phi_{v_k}(y), \quad \forall k \in N_0.$$



*Proof.* For  $s \in \mathfrak{R}^d$ , the event  $\{y_k \leq s\} \triangleq \{y_k^1 \leq s^1, \dots, y_k^d \leq s^d\}$  is considered

$$\begin{aligned}
 \bar{P}\left(\{y_k \leq S\} | \mathcal{G}_{k-1}\right) &= \bar{E}\left[I_{\{\omega: y_k(\omega) \leq s\}} | \mathcal{G}_{k-1}\right] \\
 &= \frac{E\left[\Lambda_k^{-1} I_{\{\omega: y_k(\omega) \leq s\}} | \mathcal{G}_{k-1}\right]}{E\left[\Lambda_k^{-1} | \mathcal{G}_{k-1}\right]} \\
 &= \frac{E\left[\Lambda_{k-1}^{-1} |D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} I_{\{\omega: y_k(\omega) \leq s\}} | \mathcal{G}_{k-1}\right]}{E\left[\Lambda_{k-1}^{-1} |D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} | \mathcal{G}_{k-1}\right]} \\
 &= \frac{E\left[|D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} I_{\{\omega: y_k(\omega) \leq s\}} | \mathcal{G}_{k-1}\right]}{E\left[|D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} | \mathcal{G}_{k-1}\right]}.
 \end{aligned}$$

But given

$$E\left[|D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} | \mathcal{G}_{k-1}\right] = \int_{\mathfrak{R}^d} |D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} \Phi_{v_k}(v_k) dv_k = \int_{\mathfrak{R}^d} \Phi_{v_k}(y_k) dy_k = 1$$

(using change of variables,  $dy_k = |D_k| dv_k$ ). Therefore

$$\bar{P}\left(\{y_k \leq S\} | \mathcal{G}_{k-1}\right) = \int_{\mathfrak{R}^d} I_{\{\omega: y_k(\omega) \leq s\}} \Phi_{v_k}(y_k) dy_k = \int_{-\infty}^{s^1} \dots \int_{-\infty}^{s^d} \Phi_{v_k}(y_k) dy_k, \quad q.e.d.$$

□

**Remark 2.2.7.** The above result shows the following.

- I. Start with a complete probability space  $(\Omega, \mathcal{F}, P, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\})$  on which (2.20), (2.21) are defined.

Define  $\Lambda_k^{-1} = \prod_{s=0}^k |D_s| \frac{\Phi_{v_s}(y_s)}{\Phi_{v_s}(v_s)}$ .

Then  $\Lambda_k^{-1}$  defines a new measure  $\bar{P} \ll P$  under which:

1. The distribution of  $\{x_k\}$  is the same under  $\bar{P}$  and  $P$ .
2. The sequence  $\{y_k\}$  is an  $(\{\mathcal{G}_k\}, \bar{P})$  Gaussian sequence with density

$$\Phi_{y_k}(y) = \frac{e^{-y^T y/2}}{(2\pi)^{d/2}} = \Phi_{v_k}(y), \quad \forall k \in N_0.$$

3. The inverse of  $\Lambda_k^{-1}$  can be used to define a new measure  $P \ll \bar{P}$  by setting

$$\frac{dP}{d\bar{P}} \Big|_{\mathcal{G}_k} = \Lambda_k = \prod_{s=0}^k |D_s|^{-1} \frac{\Phi_{v_s}(v_s)}{\Phi_{v_s}(y_s)}.$$

II. Start with a complete probability space  $(\Omega, \mathcal{F}, \bar{P}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\})$  on which  $\{x_k\}, \{y_k\}, k \in N_0$  are defined by

$$x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1}, \quad x_0 \in \mathbb{R}^n$$

$$\{y_k\} \text{ is an iid sequence with density } \Phi_{y_k}(y) = \frac{e^{-y^T y/2}}{(2\pi)^{d/2}}.$$

Let

$$\Lambda_k = (\Lambda_k^{-1})^{-1} = \prod_{s=0}^k |D_s|^{-1} \frac{\Phi_{v_s}(D_s^{-1}(y_s - C_s x_s))}{\Phi_{v_s}(y_s)}.$$

Then  $\frac{dP}{d\bar{P}} \Big|_{\mathcal{G}_k} = \Lambda_k$  defines a new measure  $P \ll \bar{P}$  under which:

1. The distribution of  $\{x_k\}$  is the same under  $P$  and  $\bar{P}$ .
2. The sequence  $\{v_k\}$  is an  $(\{\mathcal{G}_k\}, P)$  Gaussian sequence with density

$$\Phi_{v_k}(v) = \frac{e^{-v^T v}}{(2\pi)^{d/2}} = \Phi_{v_k}(v),$$

$$\text{where } v_k \triangleq D_k^{-1}(y_k - C_k x_k) \Rightarrow y_k = C_k x_k + D_k v_k.$$

### B. Measure Change of the State and Observation Process

Start with a complete probability space  $(\Omega, \mathcal{F}, P, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\})$  on which  $\{x_k\}, \{y_k\}, k \in N_0$  are defined by

$$x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1}, \quad x_0 \in \mathbb{R}^n, \tag{2.25}$$

$$y_k = C_k x_k + D_k v_k, \quad y_0 \in \mathbb{R}^d, \tag{2.26}$$

where

(i)  $w : \Omega \times N_0 \rightarrow \mathbb{R}^n$  is an iid sequence with density

$$\Phi_{w_k}(w) = \frac{1}{(2\pi)^{n/2}} e^{-\frac{w^T w}{2}}, \quad w_k \sim N(0, I_n);$$

(ii)  $v : \Omega \times N_0 \rightarrow \mathbb{R}^d$  is an iid sequence with density

$$\Phi_{v_k}(v) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{v^T v}{2}}, \quad v_k \sim N(0, I_d);$$

- (iii)  $x_0 : \Omega \rightarrow \mathfrak{R}^n$  a RV;
- (iv)  $\{x_0, \{w_k\}, \{v_k\}\}$  mutually independent;
- (v)  $D_k > 0, \forall k \in N_0$ ;
- (vi)  $B_k B_k^T > 0, \forall k \in N_0$ .

Next, define

$$\left. \frac{d\bar{P}}{dP} \right|_{\mathcal{G}_k} \triangleq \Lambda_k^{-1} = \prod_{s=0}^k \lambda^{-1},$$

where  $\lambda_k^{-1} = |D_k| \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)}$  and set  $\left. \frac{d\bar{P}}{dP} \right|_{\mathcal{G}_k} = \Lambda_k^{-1}$ . Then, a new probability measure is introduced

$$\left( \Omega, \mathcal{F}, \bar{P}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1} \\ \{y_k\} \text{ is an iid seq. with density } \Phi_{y_k}(y) = \Phi_{v_k}(y). \end{cases}$$

Next, start with  $(\{\mathcal{G}_k\}, \bar{P})$  and define a new measure  $\bar{\bar{P}}$  under which

$$\left( \Omega, \mathcal{F}, \bar{\bar{P}}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} \{x_k\} \text{ is an iid seq. with density} \\ \Phi_{x_k}(x) = \Phi_{w_k}(w) = (2\pi)^{-n/2} e^{-x^T x/2} \\ \{y_k\} \text{ is an iid seq. with density} \\ \Phi_{y_k}(y) = \Phi_{v_k}(y) = (2\pi)^{-d/2} e^{-y^T y/2} \end{cases}$$

To this end, the  $(\{\mathcal{G}_k\}, \bar{P})$  RV is considered, the following is defined

$$\mu_k^{-1} \triangleq |B_k| \frac{\Phi_{w_k}(x_k)}{\Phi_{w_k}(w_k)}, \quad k \in N_0$$

and set

$$M_k^{-1} = \prod_{s=1}^k |B_s| \frac{\Phi_{w_s}(x_s)}{\Phi_{w_s}(w_s)},$$

then  $\bar{E}[M_k^{-1}] = 1, \bar{P} - a.s., \forall k \in N_0$ .

Hence,  $M_k^{-1}$  can be used to define a new probability measure on  $(\Omega, \{\mathcal{G}_k\})$  through the RND

$$\left. \frac{d\bar{P}}{dP} \right|_{\mathcal{G}_k} = M_k^{-1} \quad \text{or} \quad d\bar{P}(A) = \int_A M_k^{-1} d\bar{P}(\omega), \quad \forall A \in \mathcal{G}_k.$$

Notice that

$$d\bar{P} = M_k^{-1} d\bar{P} = M_k^{-1} \frac{d\bar{P}}{dP} dP = M_k^{-1} \Lambda_k^{-1} dP$$

and hence

$$\frac{d\bar{P}}{d\bar{P}} \frac{d\bar{P}}{dP} \Big|_{\mathcal{G}_k} = \frac{d\bar{P}}{dP} \Big|_{\mathcal{G}_k} = M_k^{-1} \Lambda_k^{-1}$$

defines a change of measure from  $P$  to  $\bar{P}$ .

Next it is shown that under the probability space  $(\Omega, \mathcal{F}, \bar{P}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\})$ ,  $\{x_k\}$  is an iid seq. with density  $\Phi_{x_k}(x) = \Phi_{w_k}(w)$ , and  $\{y_k\}$  is an iid seq. with density  $\Phi_{y_k}(y) = \Phi_{v_k}(y)$ .

**Lemma 2.2.8.** *Let  $(\Omega, \mathcal{F}, \bar{P})$  be a probability space on which*

$$x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1}, \quad x_0 \in \mathfrak{R}^n$$

*$\{y_k\}$  is an iid seq. with density  $\Phi_{y_k}(y) = \Phi_{v_k}(y)$ .*

Define

$$\left. \frac{d\bar{P}}{dP} \right|_{\mathcal{G}_k} = M_k^{-1} = \prod_{s=1}^k |B_s| \frac{\Phi_{w_s}(x_s)}{\Phi_{w_s}(w_s)}.$$

Then under measure  $\bar{P}$  the seq.  $\{x_k\}$  is iid with density

$$\Phi_{x_k}(x) = \Phi_{w_k}(x) = \frac{e^{-x^T x/2}}{(2\pi)^{n/2}}, \quad \forall k \in N_0.$$

*Proof.* For  $s \in \mathfrak{R}^n$  the event  $\{x_k \leq s\}$  is considered

$$\begin{aligned} \bar{P}(\{x_k \leq s\} | \mathcal{G}_{k-1}) &= \bar{E} \left[ I_{\{\omega: x_k(\omega) \leq s\}} \Big| \mathcal{G}_{k-1} \right] \\ &= \frac{\bar{E} \left[ I_{\{\omega: x_k(\omega) \leq s\}} M_k^{-1} \Big| \mathcal{G}_{k-1} \right]}{\bar{E} \left[ M_k^{-1} \Big| \mathcal{G}_{k-1} \right]} \end{aligned}$$

$$\begin{aligned}
 &= \frac{\bar{E} \left[ M_{k-1}^{-1} |B_k| \frac{\Phi_{w_k}(x_k)}{\Phi_{w_k}(w_k)} I_{\{\omega: x_k(\omega) \leq s\}} \middle| \mathcal{G}_{k-1} \right]}{E \left[ M_{k-1}^{-1} |B_k| \frac{\Phi_{w_k}(x_k)}{\Phi_{w_k}(w_k)} \middle| \mathcal{G}_{k-1} \right]} \\
 &= \frac{\bar{E} \left[ |B_k| \frac{\Phi_{x_k}(x_k)}{\Phi_{w_k}(w_k)} I_{\{\omega: x_k(\omega) \leq s\}} \middle| \mathcal{G}_{k-1} \right]}{\bar{E} \left[ |B_k| \frac{\Phi_{w_k}(x_k)}{\Phi_{w_k}(w_k)} \middle| \mathcal{G}_{k-1} \right]}.
 \end{aligned}$$

But given

$$\begin{aligned}
 \bar{E} \left[ |B_k| \frac{\Phi_{w_k}(x_k)}{\Phi_{w_k}(w_k)} \middle| \mathcal{G}_{k-1} \right] &= \int_{\mathfrak{R}^n} |B_k| \frac{\Phi_{w_k}(x_k)}{\Phi_{w_k}(w_k)} \Phi_{w_k}(w_k) dw_k = \int_{\mathfrak{R}^n} |B_k| \Phi_{w_k}(x_k) dw_k \\
 &= \int_{\mathfrak{R}^n} \Phi_{w_k}(x_k) dx_k = 1
 \end{aligned}$$

(using the change of variables,  $w_k = B_k^{-1}(x_k - A_k x_{k-1}) \Rightarrow dw_k = |B_k|^{-1} dx_k$ ), which implies the desire result.  $\square$

**Remark 2.2.9.** *The above results show the following*

I. *Starting with*

$$\left( \Omega, \mathcal{F}, \bar{P}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1} \\ \{y_k\} \text{ is an iid seq. with density } \Phi_{y_k}(y) \sim N(0; I_d). \end{cases}$$

Let

$$M_k^{-1} = \prod_{s=1}^k |B_s| \frac{\Phi_{w_s}(x_s)}{\Phi_{w_s}(w_s)}.$$

Then  $M_k^{-1}$  defines a new measure  $\bar{\bar{P}} \ll \bar{P}$  (e.g.,  $\frac{d\bar{\bar{P}}}{d\bar{P}} \Big|_{\mathcal{G}_k} = M_k^{-1}$ ) under which:

1. *The distribution of  $\{y_k\}$  is invariant.*
2.  *$\{x_k\}$  is an  $(\mathcal{G}_k, \bar{\bar{P}})$  Gaussian seq. with density  $N(0; I_n)$ .*
3. *The inverse of  $M_k^{-1}$  can be used to define a new probability measure  $\bar{\bar{P}} \ll \bar{P}$  by setting*

$$\frac{d\bar{\bar{P}}}{d\bar{P}} \Big|_{\mathcal{G}_k} = M_k = \prod_{s=1}^k \frac{\Phi_{w_s}(B_s^{-1}(x_s - A_s x_{s-1}))}{\Phi_{w_s}(x_s)}$$

where  $w_k \triangleq B_k^{-1}(x_k - A_k x_{k-1})$  is an iid seq. with density  $N(0; I_n)$ .

II. Start with

$$\left( \Omega, \mathcal{F}, \overline{\overline{P}}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} \{x_k\} \text{ is an } \sim N(0; I_n) \text{ iid seq.} \\ \{y_k\} \text{ is an } \sim N(0; I_d) \text{ iid seq.} \end{cases}$$

Let

$$M_k = \prod_{s=1}^k \frac{\Phi_{w_s}(B_s^{-1}(x_s - A_s x_{s-1}))}{\Phi_{w_s}(x_s)}.$$

Then  $M_k$  defines a new probability measure by setting  $\left. \frac{d\overline{\overline{P}}}{dP} \right|_{\mathcal{G}_k} = M_k$  under which:

1. The distribution of  $\{y_k\}$  is invariant.
2.  $\{w_k\}$  is an  $(\mathcal{G}_k, \overline{\overline{P}})$  iid seq. with density  $N(0; I_n)$  where
 
$$w_k \triangleq B_k^{-1}(x_k - A_k x_{k-1}) \Rightarrow x_{k+1} = A_{k+1} x_k + B_{k+1} w_{k+1}.$$

### C. Change of Measure from iid Sequences

Here it is shown how the dynamics

$$\left( \Omega, \mathcal{F}, P, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} x_{k+1} = A_{k+1} x_k + B_{k+1} w_{k+1}, & x_0 \in \mathbb{R}^n \\ y_k = C_k x_k + D_k v_k & y_0 \in \mathbb{R}^d \end{cases} \quad (2.27)$$

of (2.20), (2.21) can be introduced starting with an initial probability space

$$\left( \Omega, \mathcal{F}, \overline{\overline{P}}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} x_k \sim N(0; I_n) = \Phi_{w_k}(x) = \frac{e^{-x^T x/2}}{(2\pi)^{n/2}} \\ y_k \sim N(0; I_d) = \Phi_{v_k}(y) = \frac{e^{-y^T y/2}}{(2\pi)^{d/2}}. \end{cases}$$

Define the sigma fields

$$\begin{aligned} \mathcal{G}_k &\triangleq \sigma\{x_0, \dots, x_k, y_0, \dots, y_k\} \\ \mathcal{F}_k^y &\triangleq \sigma\{y_0, \dots, y_k\}, \end{aligned}$$

where  $\mathcal{G}_k, \mathcal{F}_k^y$  are complete and the complete filtrations  $\{\mathcal{G}_k\}, \{\mathcal{F}_k^y\}$ .

Let

$$\lambda_0 = \frac{\Phi_{v_0}(D_0^{-1}(y_0 - C_0 x_0))}{|D_0| \Phi_{v_0}(y_0)}$$

and

$$\lambda_k = \frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} \times \frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{|B_k| \Psi_{w_k}(x_k)}, \quad k \geq 1.$$

Set

$$\Lambda_k = \prod_{s=0}^k \lambda_s. \quad (2.28)$$

Then define a new measure  $P$  on  $(\Omega, \{\mathcal{G}_k\})$  by setting

$$\frac{dP}{dP} \Big|_{\mathcal{G}_k} = \Lambda_k. \quad (2.29)$$

It will be shown that under  $P$ ,

$$\begin{aligned} v_k &\triangleq D_k^{-1}(y_k - C_k x_k), \quad k \in N_0 \\ w_k &\triangleq B_k^{-1}(x_k - A_k x_{k-1}), \quad k \in N_0 \end{aligned}$$

are iid seq with density  $N(0; I_d)$ ,  $(0; I_n)$ , respectively.

**Lemma 2.2.10.** *On  $(\Omega, \mathcal{F})$  and under measure  $P$  the seq. defined by*

$$\begin{aligned} v_k &\triangleq D_k^{-1}(y_k - C_k x_k), \quad k \in N_0 \\ w_k &\triangleq B_k^{-1}(x_k - A_k x_{k-1}), \quad k = 1, \dots, K \end{aligned}$$

are iid and

$$\Phi_{v_k}(v) = \frac{e^{-v^T v/2}}{(2\pi)^{d/2}}, \quad \Psi_{w_k} = \frac{e^{-w^T w/2}}{(2\pi)^{n/2}}.$$

*Proof.*  $f : \mathfrak{R}^d \rightarrow \mathfrak{R}$ ,  $g : \mathfrak{R}^n \rightarrow \mathfrak{R}$  are supposed to be bounded and Borel measurable.

It is sufficient to show

$$E \left[ g(w_k) f(v_k) \Big| \mathcal{G}_{k-1} \right] = \int_{\mathfrak{R}^d} \Phi_{v_k}(v) f(v) dv \times \int_{\mathfrak{R}^n} \Psi_{w_k}(w) g(w) dw.$$

Starting with,

$$E \left[ g(w_k) f(v_k) \Big| \mathcal{G}_{k-1} \right] = \frac{\overline{E} \left[ g(w_k) f(v_k) \Lambda_k \Big| \mathcal{G}_{k-1} \right]}{\overline{E} \left[ \Lambda_k \Big| \mathcal{G}_{k-1} \right]} = \frac{\overline{E} \left[ \lambda_k g(w_k) f(v_k) \Big| \mathcal{G}_{k-1} \right]}{\overline{E} \left[ \lambda_k \Big| \mathcal{G}_{k-1} \right]}$$

(The fact that  $\Lambda_{k-1}$  is  $\mathcal{G}_{k-1}$  measurable is used above).

Also,

$$\begin{aligned}
 \overline{\overline{E}}\left[\lambda_k \middle| \mathcal{G}_{k-1}\right] &= \overline{\overline{E}}\left[\frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} \times \frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{|B_k| \Psi_{w_k}(x_k)} \middle| \mathcal{G}_{k-1}\right] \\
 &= \overline{\overline{E}}\left\{\overline{\overline{E}}\left[\frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} \right. \right. \\
 &\quad \left. \left. \times \frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{|B_k| \Psi_{w_k}(x_k)} \middle| \mathcal{G}_{k-1}, x_k\right] \middle| \mathcal{G}_{k-1}\right\} \\
 &= \overline{\overline{E}}\left\{\frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{|B_k| \Psi_{w_k}(x_k)} \overline{\overline{E}}\right. \\
 &\quad \left. \times \left[\frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} \middle| \mathcal{G}_{k-1}, x_k\right] \middle| \mathcal{G}_{k-1}\right\}.
 \end{aligned}$$

Furthermore,

$$\overline{\overline{E}}\left[\frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} \middle| \mathcal{G}_{k-1}, x_k\right] = \frac{1}{|D_k|} \int_{\mathfrak{R}^d} \frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{\Phi_{v_k}(y_k)} \Phi_{v_k}(y_k) dy_k = 1.$$

Therefore,

$$\begin{aligned}
 \overline{\overline{E}}\left[\lambda_k \middle| \mathcal{G}_{k-1}\right] &= \frac{1}{|B_k|} \int_{\mathfrak{R}^n} \frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{\Psi_{w_k}(x_k)} \Psi_{w_k}(x_k) dx_k \\
 &= \frac{1}{|B_k|} \int_{\mathfrak{R}^n} \Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1})) dx_k = 1.
 \end{aligned}$$

Consequently,

$$\begin{aligned}
 E\left[g(w_k) f(v_k) \middle| \mathcal{G}_{k-1}\right] &= \overline{\overline{E}}\left[\lambda_k g(w_k) f(v_k) \middle| \mathcal{G}_{k-1}\right] \\
 &= \overline{\overline{E}}\left[\frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{|B_k| \Psi_{w_k}(x_k)} \frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} \right. \\
 &\quad \left. \times g(B_k^{-1}(x_k - A_k x_{k-1})) f(D_k^{-1}(y_k - C_k x_k)) \middle| \mathcal{G}_{k-1}\right] \\
 &= \overline{\overline{E}}\left\{\overline{\overline{E}}\left[\frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{|B_k| \Psi_{w_k}(x_k)} \frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} \right. \right. \\
 &\quad \left. \left. \times g(B_k^{-1}(x_k - A_k x_{k-1})) f(D_k^{-1}(y_k - C_k x_k)) \middle| \mathcal{G}_{k-1}, x_k\right] \middle| \mathcal{G}_{k-1}\right\} \\
 &= \overline{\overline{E}}\left\{\frac{\Psi_{w_k}(B_k^{-1}(x_k - A_k x_{k-1}))}{|B_k| \Psi_{w_k}(x_k)} g(B_k^{-1}(x_k - A_k x_{k-1})) \right. \\
 &\quad \left. \times \overline{\overline{E}}\left[\frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} f(D_k^{-1}(y_k - C_k x_k)) \middle| \mathcal{G}_{k-1}, x_k\right] \middle| \mathcal{G}_{k-1}\right\}.
 \end{aligned}$$



Isolating

$$\begin{aligned}
 & \overline{\overline{E}} \left[ \frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{|D_k| \Phi_{v_k}(y_k)} f(D_k^{-1}(y_k - C_k x_k)) \Big| \mathcal{G}_{k-1}, x_k \right] \\
 &= \frac{1}{|D_k|} \int_{\mathfrak{R}^d} \frac{\Phi_{v_k}(D_k^{-1}(y_k - C_k x_k))}{\Phi_{v_k}(y_k)} f(D_k^{-1}(y_k - C_k x_k)) \Phi_{v_k}(y_k) dy_k \\
 &= \frac{1}{|D_k|} \int_{\mathfrak{R}^d} \Phi_{v_k}(v_k) f(v_k) dv_k |D_k| \quad [dy_k = |D_k| dv_k]
 \end{aligned}$$

which is independent of  $x_0, x_1, \dots, x_{k-1}, y_0, \dots, y_{k-1}$  e.g., it is  $\mathcal{G}_{k-1}$  independent.

Therefore,

$$E[g(w_k) f(v_k) \Big| \mathcal{G}_{k-1}] = \int_{\mathfrak{R}^d} \Phi_{v_k}(v) f(v) dv \times \int_{\mathfrak{R}^n} \Psi_{w_k}(w) g(w) dw.$$

□

#### D. Change of Drift and Signal

The following two systems are considered

$$\left( \Omega, \mathcal{F}, P, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} x_{k+1} = A_{k+1} x_k + B_{k+1} w_{k+1}, & x_0 \in \mathfrak{R}^n \\ y_k = C_k x_k + D_k v_k, & y_0 \in \mathfrak{R}^d \end{cases} \quad (2.30)$$

$$\left( \Omega, \mathcal{F}, \tilde{P}, \{\tilde{\mathcal{G}}_k\}, \{\tilde{\mathcal{F}}_k^y\} \right) : \begin{cases} x_{k+1} = \tilde{A}_{k+1} x_k + \tilde{B}_{k+1} w_{k+1}, & x_0 \in \mathfrak{R}^n \\ y_k = \tilde{C}_k x_k + \tilde{D}_k v_k, & y_0 \in \mathfrak{R}^d. \end{cases} \quad (2.31)$$

Next, it will be shown how to define system (2.30) from (2.31), which is equivalent to determining  $\frac{dP}{d\tilde{P}} \Big|_{\mathcal{G}_k}$ .

Clearly,

$$\frac{dP}{d\tilde{P}} \Big|_{\mathcal{G}_k} = \frac{dP}{d\bar{P}} \frac{d\bar{P}}{d\tilde{P}} \Big|_{\mathcal{G}_k},$$

where  $\bar{P}$  is a measure under which  $\{x_k\}, \{y_k\}$  are iid and normal.

Thus,

$$\frac{dP}{d\tilde{P}} \Big|_{\mathcal{G}_k} = \Lambda_k = \prod_{s=0}^k \lambda_s,$$

where

$$\lambda_0 = \frac{\Phi_{v_0}(D_0^{-1}(y_0 - C_0x_0))}{|D_0|\Phi_{v_0}(y_0)}$$

and

$$\lambda_k = \frac{\Phi_{v_k}(D_k^{-1}(y_k - C_kx_k))}{|D_k|\Phi_{v_k}(y_k)} \times \frac{\Psi_{w_k}(B_k^{-1}(x_k - A_kx_{k-1}))}{|B_k|\Psi_{w_k}(x_k)}, \quad k \geq 1.$$

Also,

$$\left. \frac{d\bar{P}}{d\tilde{P}} \right|_{\mathcal{G}_k} = \tilde{\Lambda}_k = \prod_{s=0}^k \tilde{\lambda}_s,$$

where

$$\tilde{\lambda}_0 = \frac{|\tilde{D}_0|\Phi_{v_0}(y_0)}{\Phi_{v_0}(\tilde{D}_0^{-1}(y_0 - \tilde{C}_0x_0))}$$

and

$$\tilde{\lambda}_k = \frac{|\tilde{D}_k|\Phi_{v_k}(y_k)}{\Phi_{v_k}(\tilde{D}_k^{-1}(y_k - \tilde{C}_kx_k))} \times \frac{|\tilde{B}_k|\Psi_{w_k}(x_k)}{\Psi_{w_k}(\tilde{B}_k^{-1}(x_k - \tilde{A}_kx_{k-1}))}.$$

Therefore,

$$\left. \frac{dP}{d\tilde{P}} \right|_{\mathcal{G}_k} = \prod_{s=0}^k \frac{|\tilde{D}_s|\Phi_{v_s}(D_s^{-1}(y_s - C_sx_s))}{|D_s|\Phi_{v_s}(\tilde{D}_s^{-1}(y_s - \tilde{C}_sx_s))} \frac{|\tilde{B}_s|\Psi_{w_s}(B_s^{-1}(x_s - A_sx_{s-1}))}{|B_s|\Psi_{w_s}(\tilde{B}_s^{-1}(x_s - \tilde{A}_sx_{s-1}))}. \quad (2.32)$$

This is the complete data likelihood function between models (2.30) and (2.31); it is complete because it is a function of  $\mathcal{G}_k = \{x_0, \dots, x_k, y_0, \dots, y_k\}$ .

### 2.2.3 Change of Probability Measure for Nonlinear Systems

A complete probability space  $(\Omega, \mathcal{F}, P)$ , is considered, on which  $\{x_k\}, \{y_k\}$ ,  $k \in N_0$  are defined by

$$\begin{aligned} x_{k+1} &= f(x_k, w_{k+1}) \quad x_0 \in \mathfrak{R}^n \\ y_k &= h(x_k, v_k), \quad y_0 \in \mathfrak{R}^d \end{aligned} \quad (2.33)$$

where the following conditions hold:

- (i)  $x : \Omega \times N_0 \rightarrow \mathfrak{R}^n, y : \Omega \times N_0 \rightarrow \mathfrak{R}^d$ ;
- (ii)  $w : \Omega \times N_0 \rightarrow \mathfrak{R}^n$  is an indep. seq. of random variables with density  $\Psi_{w_k}(w)$ ;
- (iii)  $v : \Omega \times N_0 \rightarrow \mathfrak{R}^d$  is an indep. seq. of random variables with density  $\Phi_{v_k}(v) > 0$ ;
- (iv)  $x_0 : \Omega \rightarrow \mathfrak{R}^n$  has a density  $\Pi_0(x)$ ;
- (v)  $f : \mathfrak{R}^n \times \mathfrak{R}^n \rightarrow \mathfrak{R}^n, h : \mathfrak{R}^n \times \mathfrak{R}^d \rightarrow \mathfrak{R}^d$  are Borel measurable;
- (vi)  $\exists$  an inverse  $D : \mathfrak{R}^n \times \mathfrak{R}^n \rightarrow \mathfrak{R}^n$  such that  $w_{k+1} = D(x_{k+1}, x_k)$
- (vii)  $\exists$  an inverse  $G : \mathfrak{R}^d \times \mathfrak{R}^n \rightarrow \mathfrak{R}^n$  such that  $v_k = G(y_k, x_k)$ ;
- (viii) the derivatives  $\left. \frac{\partial}{\partial y} G(y, x_e) \right|_{y=y_e}, \left. \frac{\partial}{\partial v} h(x_e, v) \right|_{v=v_e}$  are continuous and nonsingular.

Let

$$\begin{aligned} \mathcal{G}_k &\triangleq \sigma\{x_0, x_1, \dots, x_k, y_0, \dots, y_k\} \\ \mathcal{F}_k^y &\triangleq \sigma\{y_0, \dots, y_k\} \end{aligned}$$

be complete  $\sigma$ -algebras.

Define

$$\Lambda_k^{-1} = \prod_{s=0}^k \frac{\Phi_{v_s}(y_s)}{\Phi_{v_s}(v_s)} \left. \frac{\partial}{\partial y} G(y, x_s) \right|_{y=y_s}^{-1}.$$

Then

$$E[\Lambda_k^{-1}] = E[E[\Lambda_k^{-1} | \mathcal{G}_{k-1}]] = E[\Lambda_k^{-1} E[\frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} \left. \frac{\partial}{\partial y} G(y, x_k) \right|_{y=y_k}^{-1} | \mathcal{G}_{k-1}]].$$

But

$$\begin{aligned} E\left[\frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} \left. \frac{\partial}{\partial y} G(y, x_k) \right|_{y=y_k}^{-1} \Big| \mathcal{G}_{k-1}\right] &= \int_{\mathfrak{R}^d} \frac{\Phi_{v_k}(y_k)}{\Phi_{v_k}(v_k)} \left. \frac{\partial}{\partial y} G(y, x_k) \right|_{y=y_k}^{-1} \Phi_{v_k}(v_k) dv_k \\ &= \int_{\mathfrak{R}^d} \Phi_{v_k}(y_k) dy_k = 1 \end{aligned}$$

$$\Rightarrow E[\Lambda_k^{-1}] = E[\Lambda_{k-1}^{-1}], \quad \forall k \in N_0$$

$$\Rightarrow E[\Lambda_k^{-1}] = E[\Lambda_0^{-1}] = \int_{\mathfrak{R}^d} \frac{\Phi_{v_0}(y_0)}{\Phi_{v_0}(v_0)} \left. \frac{\partial}{\partial y} G(y, x_0) \right|_{y=y_0}^{-1} \Phi_{v_0}(v_0) dv_0 = 1.$$

Therefore, the new measure can be defined by letting

$$\frac{d\bar{P}}{dP}\Big|_{\mathcal{G}_k} = \Lambda_k^{-1}.$$

Such that under measure  $\bar{P}$ , the seq.  $\{y_k\}$  is independent with density  $\phi_{y_k}(y_k)$ .

Also, starting with  $\bar{P}$ ,  $P$  can be constructed by setting

$$\begin{aligned} \Lambda_k &= \prod_{s=0}^k \frac{\Phi_{v_s}(v_s)}{\Phi_{v_s}(y_s)} \Big|_{\frac{\partial}{\partial v} h(x_s, v)} \Big|_{v=v_s}^{-1} \\ v_s &\triangleq G(y_s, x_s) \end{aligned}$$

and defining

$$\frac{dP}{d\bar{P}}\Big|_{\mathcal{G}_k} \triangleq \Lambda_k.$$

**Lemma 2.2.11.** *The following system is considered,*

$$\left( \Omega, \mathcal{F}, \bar{P}, \{\mathcal{G}_k\}, \{\mathcal{F}_k^{\mathcal{Y}}\} \right) : \begin{cases} x_{k+1} = f(x_k, w_{k+1}), & x_0 \in \mathfrak{R}^n \\ \{y_k\} \text{ is an indep. seq. with density } \Phi_{y_k}(y_k). \end{cases}$$

Then under measure  $P$ , the seq.  $\{v_k\}$  is independent having densities  $\{\Phi_{v_k}(\cdot)\}$  and

$$\begin{aligned} v_k &\triangleq G(y_k, x_k) \\ \Rightarrow y_k &= h(x_k, v_k). \end{aligned}$$

*Proof.* Starting with

$$\begin{aligned} P(v_k \leq s | \mathcal{G}_{k-1}) &= E \left[ I_{\{\omega: v_k(\omega) \leq s\}} | \mathcal{G}_{k-1} \right] \\ &= \frac{\bar{E} \left[ I_{\{\omega: x_k(\omega) \leq s\}} \Lambda_k | \mathcal{G}_{k-1} \right]}{\bar{E} \left[ \Lambda_k | \mathcal{G}_{k-1} \right]} \\ &= \bar{E} \left[ \lambda_k I_{\{\omega: x_k(\omega) \leq s\}} | \mathcal{G}_{k-1} \right] \\ &= \int_{\mathfrak{R}^d} I_{\{\omega: x_k(\omega) \leq s\}} \frac{\Phi_{v_k}(v_k)}{\Phi_{v_k}(y_k)} \Big|_{\frac{\partial h}{\partial v}}^{-1} \Phi_{v_k}(y_k) dy \\ &= P(v_k \leq s) \end{aligned}$$

$\Rightarrow$  under  $P$ ,  $\{v_k\}$  is an indep. seq. with density  $\Phi_{v_k}$ ,  $v_k \triangleq G(y_k, x_k)$ . □

**A. Relation Between RND's and Sample Densities of RP's**

The following system is considered

$$\left( \Omega, \mathcal{F}, P, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\} \right) : \begin{cases} x_{k+1} = f(x_k) + g(x_k)w_{k+1}, & x_0 \in \mathbb{R}^n \\ y_k = h(x_k) + D_k v_k, & y_0 \in \mathbb{R}^d \end{cases}$$

with the assumptions that  $\{w_k\}, \{v_k\}$  are iid seq. with densities

$$\begin{aligned} \Phi_{w_k}(w) &= \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{w^T w}{2}}, \\ \Phi_{v_k}(v) &= \frac{1}{(2\pi)^{\frac{d}{2}}} e^{-\frac{v^T v}{2}}. \end{aligned}$$

If someone is interested in the density of the sample path  $\{y_0, \dots, y_k\}$  given the data  $\{x_0, \dots, x_k\}$ .

Then

$$\begin{aligned} P(y_0 \leq y | x_0) &= P(h(x_0) + D_0 v_0 \leq y | x_0) \\ &= P(v_0 \leq D_0^{-1}(y - h(x_0)) | x_0) \\ &= P(v_0 \leq D_0^{-1}(y - h(x_0))) \\ &= \int_{-\infty}^{D_0^{-1}(y - h(x_0))} \Phi_{v_0}(v) dv \\ \Rightarrow \frac{d}{dy} P(y_0 \leq y | x_0) &= \frac{d}{dy} \int_{-\infty}^{D_0^{-1}(y - h(x_0))} \Phi_{v_0}(v) dv \\ &= \Phi_{v_0}(D_0^{-1}(y - h(x_0))) |D_0|^{-1} \end{aligned}$$

e.g.,

$$p(y_0 | x_0) = \frac{d}{dy_0} P(y_0 | x_0) = |D_0|^{-1} \Phi_{v_0}(D_0^{-1}(y_0 - h(x_0))).$$

Similarly,

$$\begin{aligned} p(y_0, \dots, y_k | x_0, \dots, x_k) &= \frac{p(y_0, \dots, y_k, x_0, \dots, x_k)}{p(x_0, \dots, x_k)} \\ &= \frac{\prod_{s=0}^k p(y_s | y_0, \dots, y_{s-1}, x_0, \dots, x_s) p(x_s | x_0, \dots, x_{s-1}, y_0, \dots, y_{s-1})}{\prod_{s=0}^k p(x_s | x_0, \dots, x_{s-1})}. \end{aligned}$$

Given that  $x_s$  is independent of  $\{y_0, \dots, y_{s-1}\}$  and  $\{x_0, \dots, x_{s-2}\}$  given  $x_{s-1}$ , and  $y_s$  is independent of  $\{y_0, \dots, y_{s-1}\}$  and  $\{x_0, \dots, x_{s-1}\}$  given  $x_s$ , then

$$\begin{aligned} p\left(y_0, \dots, y_k \mid x_0, \dots, x_k\right) &= \frac{\prod_{s=0}^k p\left(y_s \mid x_s\right) p\left(x_s \mid x_{s-1}\right)}{\prod_{s=0}^k p\left(x_s \mid x_{s-1}\right)} \\ &= \prod_{s=0}^k P\left(y_s \mid x_s\right) \\ \Rightarrow p\left(y_0, \dots, y_k \mid x_0, \dots, x_k\right) &= \prod_{s=0}^k |D_s|^{-1} \Phi_{v_s}\left(D_s^{-1}\left(y_s - h\left(x_s\right)\right)\right). \end{aligned}$$

Therefore,

$$\begin{aligned} \frac{dP}{d\bar{P}} \Big|_{\mathcal{G}_k} = \Lambda_k &= \frac{p\left(y_0, \dots, y_k \mid x_0, \dots, x_k\right)}{p_v\left(y_0, \dots, y_k\right)} \\ &= \prod_{s=0}^k \frac{|D_s|^{-1} \Phi_{v_s}\left(D_s^{-1}\left(y_s - h\left(x_s\right)\right)\right)}{\Phi_{v_s}\left(y_s\right)} \end{aligned}$$

defines a measure  $P$ , starting with measure  $\bar{P}$  under which  $\{y_k\}$  is an iid seq. with density  $\Phi_{v_k}(v)$ .

Similarly, suppose the sample path data  $\{y_0, \dots, y_k, x_0, \dots, x_k\}$  are given. Then its sample path density is

$$p\left(y_0, \dots, y_k, x_0, \dots, x_k\right) = p\left(y_0, \dots, y_k \mid x_0, \dots, x_k\right) P\left(x_0, \dots, x_k\right).$$

But

$$\begin{aligned} p\left(x_0, \dots, x_k\right) &= p\left(x_k \mid x_0, \dots, x_{k-1}\right) p\left(x_0, \dots, x_{k-1}\right) \\ &= p\left(x_0\right) \prod_{s=1}^k p\left(x_s \mid x_{s-1}\right), \end{aligned}$$

$$\begin{aligned} P\left(x_k \leq x \mid x_{k-1}\right) &= P\left(f\left(x_{k-1}\right) + g\left(x_{k-1}\right) w_k \leq x \mid x_{k-1}\right) \\ &= P\left(w_k \leq g^{-1}\left(x_{k-1}\right)\left(x - f\left(x_{k-1}\right)\right) \mid x_{k-1}\right) \\ &= \int_{-\infty}^{g^{-1}\left(x_{k-1}\right)\left(x - f\left(x_{k-1}\right)\right)} \Phi_{w_k}(w) dw. \end{aligned}$$

Let  $\frac{d}{dx}P(x_k \leq x | x_{k-1}) = p(x_k | x_{k-1})$ , then

$$\begin{aligned} p(x_k | x_{k-1}) &= |g(x_{k-1})|^{-1} \Phi_{w_k} \left( g^{-1}(x_{k-1})(x - f(x_{k-1})) \right) \\ \Rightarrow p(x_0, \dots, x_k) &= p(x_0) \prod_{s=1}^k |g(x_{s-1})|^{-1} \Phi_{w_s} \left( g^{-1}(x_{s-1})(x_s - f(x_{s-1})) \right). \end{aligned}$$

Therefore,

$$\begin{aligned} \frac{dP}{d\widehat{P}} \Big|_{\mathcal{F}_k^x} &= M_k = \frac{p(x_0, \dots, x_k)}{p_{w_k}(x_0, \dots, x_k)} \\ &= \prod_{s=1}^k \frac{|g(x_{s-1})|^{-1} \Phi_{w_s} \left( g^{-1}(x_{s-1})(x_s - f(x_{s-1})) \right)}{\Psi_{w_s}(x_s)} \end{aligned}$$

defines measure  $P$ , starting with measure  $\widehat{P}$  under which  $\{x_k\}$  is an iid seq with density  $\{\Phi_{w_k}(x_k)\}$ .

Finally, the sample density of  $\mathcal{G}_k$  given by

$$\begin{aligned} \frac{p(y_0, \dots, y_k, x_0, \dots, x_k)}{p_{w_k}(x_0, \dots, x_k) p_{v_k}(y_0, \dots, y_k)} &= \prod_{s=0}^k \frac{|g(x_{s-1})|^{-1} \Phi_{w_s} \left( g^{-1}(x_{s-1})(x_s - f(x_{s-1})) \right)}{\Psi_{w_s}(x_s)} \\ &\quad \times \frac{|D_s|^{-1} \Phi_{v_s} \left( D_s^{-1}(y_s - h(x_s)) \right)}{\Phi_{v_s}(y_s)} \end{aligned}$$

defines a change of measure from  $\overline{\overline{P}}$  under which  $\{y_k\}, \{x_k\}$  are iid with densities  $\Phi_{v_k}(\cdot), \Psi_{w_k}(\cdot)$ , respectively, to measure  $P$ .

### 2.2.4 Nonlinear Filtering Prediction and Smoothing

Here, the nonlinear filtering, smoothing and prediction are presented. It is assumed that there is a signal  $\{x_k\}$ , called the state of the system which is not directly observable. Rather only some noisy function  $\{y_k\}$  of  $\{x_k\}$  are observed, called the observation process. The objective is to obtain an expression for the "best estimate" of  $x_s$  (or  $\Phi(x_s)$  for  $\Phi$  in a certain class of functions), given the history of the observations  $\{y_\tau : 0 \leq \tau \leq k\}$ , that is, given the observable  $\sigma$ -field.

$$\mathcal{Y}_k = \sigma\{y_\tau : 0 \leq \tau \leq k\}.$$

Depending of the relation between the times "s" and "k" the following situations exist:

1. **The filtering problem.** Given the observation  $\mathcal{Y}_k$  find the best estimate  $\hat{x}_k$  of the state  $x_k$  from these observations.
2. **The smoothing problem.** Given the observation  $\mathcal{Y}_k$  find the best estimate  $\hat{x}_s$  of the state  $x_s$ ,  $s < k$  from these observations.
3. **The prediction problem.** Given the observation  $\mathcal{Y}_k$  find the best estimate  $\hat{x}_s$  of the state  $x_s$ ,  $s > k$  from these observations.

The computations should be done recursively, in terms of a static  $\{\pi_k\}$  which can be updated using only new observations

$$\pi_{s+\tau} = \alpha(s, \tau, \pi_s, \{y_{k+u} : 0 \leq u \leq \tau\}) \quad (2.34)$$

in which  $s = k$  corresponds to the filtering problem,  $s > k$  corresponds to the prediction problem, and  $s < k$  corresponds to the smoothing problem. The statistic  $\{\pi_k\}$  is then used to calculate, pointwise estimates  $\hat{\Phi}(x_s)$  of functionals  $\Phi(x_s)$  from the observations  $\mathcal{Y}_k$ :

$$\hat{\Phi}(x_s) = \beta(k, s, \Phi, y_k, \pi_s). \quad (2.35)$$

Suppose the processes  $\{x_k\}$ ,  $\{y_k\}$  are defined on a fixed probability space  $(\Omega, \mathcal{F}, P)$  with filtration  $\{\mathcal{F}_k\}$  and finite time  $k \in N_0^m \triangleq \{0, 1, \dots, m\}$ . The specification of the best estimate  $\hat{x}_s$  of  $x_s$  from  $\mathcal{Y}_k$  is usually done by minimizing a functional of the distance of  $x_s$  from the closed subspace generated by  $\mathcal{Y}_k \triangleq \{y_s : 0 \leq s \leq k\}$ . From the general theory of Hilbert spaces it is known that  $\hat{x}_s(\mathcal{Y}_k)$  is the orthogonal projection onto the subspace generated by  $\mathcal{Y}_k$ .

**Theorem 2.2.12.** *Let  $(\Omega, \mathcal{F}, P)$  be a probability space with filtration  $\{\mathcal{F}_k\}$ , a sub-sigma algebra  $\{\mathcal{F}_k^y\}$  of  $\{\mathcal{F}_k\}$  and  $X \in L^2(\Omega, \mathcal{F}, P)$ .*

Define

$$M_{\mathcal{Y}_k} \triangleq \{y(j, \cdot) : \Omega \rightarrow \mathfrak{R}^d : y(j, \cdot) \in L^2(\Omega, \mathcal{F}, P) \text{ and } y_j \text{ is } \mathcal{F}_j^y - \text{measurable, } j \in N_0^k\},$$

$$\mathcal{P}_{M_{\mathcal{Y}_k}} : L^2(\Omega, \mathcal{F}, P) \rightarrow M_{\mathcal{Y}_k}.$$



Then

$$\mathcal{P}_{M_{\mathcal{Y}_k}}(x_s) = E[x_s | M_{\mathcal{Y}_k}] = E[x_s | \mathcal{F}_k^{\mathcal{Y}}]$$

is  $P$ -a.s. unique, e.g.,

$$\int_A \mathcal{P}_{M_{\mathcal{Y}_k}}(x_s)(\omega) dP(\omega) = \int_A x_s(\omega) dP(\omega), \quad \forall A \in \mathcal{F}_k^{\mathcal{Y}}.$$

**Remark 2.2.13.** The minimum variance estimate  $\mathcal{P}_{M_{\mathcal{Y}_k}}(x_s)$  is unbiased, that is,  $E[\mathcal{P}_{M_{\mathcal{Y}_k}}(x_s)] = E[x_s]$ ,  $\forall k \in N_0$ .

**Theorem 2.2.14.** (The minimum variance estimation) Given the observations  $\mathcal{Y}_k$ , the minimum variance estimate of the state  $x_s$  defined by

$$\inf \left\{ E \left[ \|x_s - m_k\|^2 : m_k \in M_{\mathcal{Y}_k} \right] \right\} \quad (2.36)$$

is given by  $\hat{x}_s(y_k) = E[x_s | \mathcal{F}_k^{\mathcal{Y}}]$ .

*Proof.* Since

$$E \left[ \|x_s - m_k\|^2 \right] = E \left[ \|x_s - \hat{x}_s(\mathcal{Y}_k) + \hat{x}_s(\mathcal{Y}_k) - m_k\|^2 \right] \quad (2.37)$$

$$= E \left[ \|m_k - \hat{x}_s(\mathcal{Y}_k)\|^2 \right] + E \left[ \|x_s - \hat{x}_s(\mathcal{Y}_k)\|^2 \right] + 2E \left[ (x_s - \hat{x}_s(\mathcal{Y}_k), \hat{x}_s(\mathcal{Y}_k) - m_k) \right]. \quad (2.38)$$

Reconditioning on  $\mathcal{F}_k^{\mathcal{Y}}$ , the last term in (2.38) vanishes; since the second right-side term of (2.38) is fixed the left-side is minimized by setting  $m_k = \hat{x}_s(\mathcal{Y}_k)$ . Therefore, the conditional expectation  $E[x_s | \mathcal{F}_k^{\mathcal{Y}}]$  gives the least minimum variance estimate.  $\square$

**Definition 2.2.15.** (The risk-sensitive estimation) Given the observations  $\mathcal{Y}_k$ , the minimum risk-sensitive estimate  $\hat{x}_s(\mathcal{Y}_k)$  of the state  $x_s$  is defined by

$$\begin{aligned} & \inf \left\{ E \left[ \exp \left( \theta \sum_{\tau=0}^s \|x_\tau - m_\tau\|^2 : m_k \in M_{[0,k]}, \theta \in \mathfrak{R} \right) \right] \right\} \\ & = \left\{ E \left[ \exp \left( \theta \mathcal{C}_{0,k}(x_0^s, m_0^s) \right) \right] : m_k \in M_{[0,k]}, \theta \in \mathfrak{R} \right\} \end{aligned} \quad (2.39)$$

where

$$\begin{aligned} M_{[0,k]} = \left\{ y : N_0^k \times \Omega \rightarrow \mathfrak{R}^d : y(k, \omega) \in L^2((0, k) \times \Omega, \mathcal{F}, P), \text{ for almost all } k \right. \\ \left. y(k, \cdot) \in L^2(\Omega, \mathcal{F}, P) \text{ and } y_k \text{ is } \mathcal{Y}_k\text{-measurable} \right\}. \end{aligned} \quad (2.40)$$

The parameter  $\theta$  is the so-called sensitive parameter which renders the optimizer optimistic if  $\theta < 0$  and pessimistic if  $\theta > 0$ . In addition, for  $\theta > 0$  the sample cost  $\exp(\theta\mathcal{C})$  is a convex increase function of  $\mathcal{C}$ . In the former case the optimizer is concerned with designing under worst scenario, hence the name risk-averse, while in the latter case is concerned with designing under favorable conditions, hence the name risk-seeking. The case  $\theta = 0$  lies between the risk-averse and the risk-seeking optimalities.

The following expansion is considered

$$E\left[\exp(\theta\mathcal{C})\right] = 1 + \theta E\left[\mathcal{C}\right] + \frac{1}{2}\theta^2 E\left[\mathcal{C}^2\right] + O(\theta^3). \quad (2.41)$$

Then

$$\theta \log E\left[\exp(\theta\mathcal{C})\right] = E\left[\mathcal{C}\right] + \frac{1}{2}\theta \text{Var}\left(\mathcal{C}\right) + O(\theta^2). \quad (2.42)$$

Consequently, in the limit, as  $\theta \rightarrow 0$ , the risk-sensitive estimation problem includes as a special case the risk-neutral case.

**Definition 2.2.16.** Let  $(\Omega, \mathcal{F}, P, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\})$  is a complete probability space on which the state process  $\{x_k\}, k \in N_0$  and the observation process  $\{y_k\}, k \in N_0$ , are defined by the following recursions.

$$\begin{aligned} x_{k+1} &= f(k+1, x_k) + g(k+1, x_k)w_{k+1}, \quad x_0 \in \mathfrak{R}^n, \\ y_k &= h(k, x_k) + D_k v_k \quad y_0 \in \mathfrak{R}^d, \end{aligned} \quad (2.43)$$

in which the condition of (2.33) hold, and additionally

$$\begin{aligned} \Phi_k(v) &\triangleq \Phi_{v_k}(v) = \frac{1}{(2\pi)^{d/2}} e^{v^T v/2}, \\ \Psi_k(w) &\triangleq \Psi_{w_k}(w) = \frac{1}{(2\pi)^{n/2}} e^{w^T w/2}, \\ P(x_0 \leq x) &= \int_{-\infty}^x \Pi_{x_0}(z) dz. \end{aligned}$$

Clearly, the minimum variance estimate of any Borel bounded function  $\Phi : \mathfrak{R}^n \rightarrow \mathfrak{R}$  is given by

$$\Pi_k(\Phi) = E\left[\Phi(x_k) \middle| \mathcal{F}_k^y\right]. \quad (2.44)$$

Thus, a program will be introduced which will enable the computation of the conditional distribution of  $x_k$  given  $\mathcal{F}_k^y$ , namely

$$P(x_k \in A | \mathcal{F}_k^y), \quad \forall A \in \mathcal{B}(\mathfrak{R}^n). \quad (2.45)$$

**Definition 2.2.17.** For any given Borel measurable function which is bounded,  $\Phi : N_0 \times \mathfrak{R}^n \rightarrow \mathfrak{R}$ , the measure-value process  $\Pi_k(\Phi)$  is related to another measure-valued process through

$$\Pi_k(\Phi) = \frac{\overline{E}[\Phi(k, x_k)\Lambda_k | \mathcal{F}_k^{\mathcal{Y}}]}{\overline{E}[\Lambda_k | \mathcal{F}_k^{\mathcal{Y}}]} \stackrel{\nabla}{=} \frac{\overline{\Pi}_k(\Phi)}{\overline{\Pi}_k(1)}, \quad P - a.s.$$

where

$$\frac{dP}{d\overline{P}} \Big|_{\mathcal{G}_k} = \Lambda_k.$$

Let  $BC(\mathfrak{R}^n) \triangleq \{\Phi : \mathfrak{R}^n \rightarrow \mathfrak{R} : \Phi \text{ is bounded and continuous}\}$ , equipped with the norm topology  $\|\Phi\|_{\infty} \triangleq \sup\{|\Phi| : x \in \mathfrak{R}^n, \Phi \in BC(\mathfrak{R}^n)\}$ . Let  $\mathcal{M}_+(\mathfrak{R}^n) \triangleq \{\}$  denotes the set of positive measures on  $(\mathfrak{R}^n, \mathcal{B}(\mathfrak{R}^n))$ , equipped with the norm topology  $\|\mu\| = Var_A \mu = \sup_{A_i \subset A \in \mathcal{B}(\mathfrak{R}^n)} \sum_{i=1}^n \mu(A_i)$ , when the upper bound is taken over all finite collections of pairwise disjoint sets  $\{A_i\}_{i=1}^n \in \mathcal{B}(\mathfrak{R}^n)$  which are contained in  $A$ . It is noted that  $Var_A \mu = \mu^+(A) + \mu^-(A)$ ,  $\forall A \in \mathcal{B}(\mathfrak{R}^n)$ . Consequently,  $Var_A \mu$  is a measure with respect to  $A$ , which is denoted by  $|\mu|(A) = Var_A \mu$ .

If  $\mu$  is a countably additive function on the  $\sigma$ -field  $\mathcal{B}(\mathfrak{R}^n)$ , then  $\mu$  is finite if and only if:

- $\mu^+, \mu^-$  are finite,
- $|\mu(A)| < \infty$ ,
- $Var_A \mu < \infty$ .

From this, it follows that if  $|\mu|(A) < \infty$  for every  $\mathcal{B}(\mathfrak{R}^n)$  measurable function  $\Phi$  which is bounded on  $A$ , and is  $\mu$ -measurable on  $A$ , e.g.,

$$\int_A \Phi(x) d\mu(x) = \int_A \Phi(x) d\mu^+(x) - \int_A \Phi(x) d\mu^-(x)$$

and

$$\left| \int_A \Phi(x) d\mu(x) \right| \leq \sup_{x \in A} |\Phi(x)| |\mu|(A).$$

Thus, for any  $\mu \in \mathcal{M}_+(\mathfrak{R}^n)$ ,  $\Phi \in BC(\mathfrak{R}^n)$  the inner product is defined by

$$(\Phi, \mu) = \int_{\mathfrak{R}^n} \Phi(x) d\mu(x).$$

**Remark 2.2.18.** *The Bayes formula of Definition (2.2.17) is the starting point of non-linear minimum variance filtering. The process*

$$\pi_k(\Phi) \triangleq \overline{E} \left[ \Phi(x_k) \Lambda_k \middle| \mathcal{F}_k^{\mathcal{Y}} \right]$$

is the unnormalized measure-valued process of  $\Pi_k(\Phi)$  since

$$\Pi_k(\Phi) = \frac{\pi_k(\Phi)}{\pi_k(1)} = \frac{\int_{\mathbb{R}^n} \Phi(z) d\pi_k(z)}{\int_{\mathbb{R}^n} 1 d\pi_k(z)} = \frac{(\Phi, \pi)}{(1, \pi)}.$$

Further if the measure  $\pi_k(\cdot)$  is absolutely continues with respect to the Lebeggue measure, then  $\frac{d}{dx} \pi_k(x)$  exist and

$$\frac{d}{dx} \Pi_k(x) = \frac{\pi_k(x)}{\int_{\mathbb{R}^n} \pi_k(x) dx}.$$

Next a recursive equation for  $\pi_k(x)$  is derived.

The system of Definition 2.2.16, is considered, starting with measure  $\overline{P}$  under which  $\{x_k\}, \{y_k\}$  are white noise sequences.

Let

$$\begin{aligned} \mathcal{G}_k &= \sigma\{x_0, \dots, x_k, y_0, \dots, y_k\}, \quad k \in N_0 \\ \mathcal{F}_k^{\mathcal{Y}} &= \sigma\{y_0, \dots, y_k\}, \quad k \in N_0 \end{aligned}$$

where  $\{\mathcal{G}_k\}, \{\mathcal{F}_k^{\mathcal{Y}}\}, k \in N_0$  are their corresponding filtrations.

Let

$$\begin{aligned} \lambda_0 &= \frac{\Phi_{v_0} \left( D_0^{-1}(y_0 - h(0, x_0)) \right)}{|D_0| \Phi_{v_0}(y_0)}, \\ \lambda_s &= \frac{\Phi_{v_s} \left( D_s^{-1}(y_s - h(s, x_s)) \right) \Psi_{w_s} \left( g^{-1}(s, x_{s-1})(x_s - f(s, x_{s-1})) \right)}{|D_s| \Phi_{v_s}(y_s) |g(s, x_{s-1})| \Psi_{w_s}(x_s)}, \quad s \in N_0. \end{aligned}$$

Next, a new probability measure  $P$  is defined on  $(\Omega, \mathcal{F}, \{\mathcal{G}_k\})$  by introducing the RND

$$\frac{dP}{d\overline{P}} \Big|_{\mathcal{G}_k} = \Lambda_k = \prod_{s=0}^k \lambda_s.$$

Define

$$\begin{aligned} w_s &\triangleq g^{-1}(s, x_{s-1})(x_s - f(s, x_{s-1})), \quad s = 1, 2, \dots \\ v_s &\triangleq D_s^{-1}(y_s - h(s, x_s)), \quad s = 0, 1, \dots \end{aligned}$$

Then under measure  $P$ , the sequences  $\{v_k\}, \{w_k\}$  are independent,  $N(0; I_d), N(0; I_n)$ , respectively.

**Theorem 2.2.19.** (Recursive Equation) Consider  $(\Omega, \mathcal{F}, P, \{\mathcal{G}_k\}, \{\mathcal{F}_k^y\})$  a complete probability space on which  $\{x_k, y_k\}, k \in N_0$  are defined by Definition 2.2.16.

Then for any  $\Phi \in BC(\mathbb{R}^n)$  the measure valued process  $\{\pi_k(\Phi)\}, k \in N_0$  satisfies the following recursion

$$\pi_k(\Phi) = \int_{\mathbb{R}^n} \Phi(x) \frac{\Phi_k(D_k^{-1}(y_k - h(k, x)))}{|D_k|\Phi_k(y_k)} \pi_{k-1}\left(\frac{\Psi_k(g^{-1}(k, \cdot))(x - f(k, \cdot))}{|g(k, \cdot)|}\right) dx \quad (2.46)$$

with initial condition

$$\pi_0(\Phi) = \int_{\mathbb{R}^n} \Phi(x) \frac{\Phi_0(D_0^{-1}(y_0 - h(0, x)))}{|D_0|\Phi_0(y_0)} \Pi_{x_0}(x) dx. \quad (2.47)$$

Further, if  $\pi_k(\Phi)$  has a density

$$\pi_k(\Phi) = \int_{\mathbb{R}^n} \Phi(x) d\pi_k(x) = \int_{\mathbb{R}^n} \phi(x) \pi_k(x) dx$$

then  $\{\pi_k(x)\}, k \in N_0$  satisfies the recursion

$$\pi_k(x) = \frac{\Phi_k(D_k^{-1}(y_k - h(k, x)))}{|D_k|\Phi_k(y_k)} \int_{\mathbb{R}^n} \frac{\Psi_k(g^{-1}(k, z))(x - f(k, z))}{|g(k, z)|} \pi_{k-1}(z) dz \quad (2.48)$$

with initial condition

$$\pi_0(x) = \frac{1}{|D_0|\Phi_0(y_0)} \Phi_0(D_0^{-1}(y_0 - h(0, x))) \Pi_{x_0}(x). \quad (2.49)$$

*Proof.* By definition

$$\begin{aligned} \pi_k(\Phi) &= \overline{E} \left[ \Phi(x_k) \Lambda_k \middle| \mathcal{F}_k^y \right] \\ &= \int_{\mathbb{R}^n} \Phi(x) d\pi_k(x). \end{aligned} \quad (2.50)$$

Further,

$$\pi_k(\Phi) = \bar{E} \left[ \Phi(x_k) \Lambda_k \middle| \mathcal{F}_k^{\mathcal{Y}} \right] = \bar{E} \left[ \prod_{s=0}^k \Phi(x_s) \Lambda_s \middle| \mathcal{F}_k^{\mathcal{Y}} \right] = E^Q \left[ \prod_{s=0}^k \Phi(x_s) \Lambda_s \right]$$

where  $Q$  is the measure induced by the independent sequence  $\{x_s\}$  (e.g.,  $\{x_s\}$  and  $\{y_s\}$  are independent).

$$\begin{aligned} \pi_k(\Phi) &= E^Q \left[ \Lambda_{k-1} \lambda_k \Phi(x_k) \right] \\ &= E^Q \left[ \Lambda_{k-1} \frac{\Phi_k(D_k^{-1}(y_k - h(k, x_k)))}{|D_k| \Phi_k(y_k)} \frac{\Psi_k(g^{-1}(k, x_{k-1})(x_k - f(k, x_{k-1})))}{|g(k, x_{k-1})| \Psi_k(x_k)} \Phi(x_k) \right] \\ &= \frac{1}{|D_k| \Phi_k(y_k)} E^Q \left[ \Lambda_{k-1} \Phi_k(D_k^{-1}(y_k - h(k, x_k))) \right. \\ &\quad \left. \times \frac{\Psi_k(g^{-1}(k, x_{k-1})(x_k - f(k, x_{k-1})))}{|g(k, x_{k-1})| \Psi_k(x_k)} \Phi(x_k) \right] \\ &= \frac{1}{|D_k| \Phi_k(y_k)} E^Q \left[ \frac{\Lambda_{k-1}}{|g(k, x_{k-1})|} \right. \\ &\quad \left. \times E^Q \left[ \frac{\Phi_k(D_k^{-1}(y_k - h(k, x_k))) \Psi_k(g^{-1}(k, x_{k-1})(x_k - f(k, x_{k-1})))}{\Psi_k(x_k)} \right. \right. \\ &\quad \left. \left. \Phi(x_k) \middle| \mathcal{F}_{k-1}^{\mathcal{Y}}, x_{k-1} \right] \right] \\ &= \frac{1}{|D_k| \Phi_k(y_k)} E^Q \left[ \frac{\Lambda_{k-1}}{|g(k, x_{k-1})|} \right. \\ &\quad \left. \times E^Q \left[ \frac{\Phi_k(D_k^{-1}(y_k - h(k, x_k))) \Psi_k(g^{-1}(k, x_{k-1})(x_k - f(k, x_{k-1})))}{\Psi_k(x_k)} \right. \right. \\ &\quad \left. \left. \Phi(x_k) \right] \right] \\ &= \frac{1}{|D_k| \Phi_k(y_k)} E^Q \left[ \frac{\Lambda_{k-1}}{|g(k, x_{k-1})|} \right. \\ &\quad \left. \times \int_{\mathfrak{R}^n} \frac{\Phi_k(D_k^{-1}(y_k - h(k, \xi))) \Psi_k(g^{-1}(k, x_{k-1})(\xi - f(k, x_{k-1})))}{\Psi_k(\xi)} \right. \\ &\quad \left. \Phi(\xi) \Psi_k(\xi) d\xi \right] \\ &= \frac{1}{|D_k| \Phi_k(y_k)} E^Q \left[ \frac{\Lambda_{k-1}}{|g(t, x_{k-1})|} \tilde{g}(k, x_{k-1}, y_k) \right]. \end{aligned}$$

Define

$$\begin{aligned} \tilde{g}(k, x_{k-1}, y_k) &\triangleq \\ &\int_{\mathfrak{R}^n} \Phi_k(D_k^{-1}(y_k - h(k, \xi))) \Psi_{w_k}(g^{-1}(k, x_{k-1})(\xi - f(k, x_{k-1}))) \Phi(\xi) d\xi. \end{aligned}$$

Then

$$\pi_k(\Phi) = \frac{1}{|D_k|\Phi_k(y_k)} \int_{\mathbb{R}^n} \frac{\tilde{g}(k, z, y_k)}{|g(k, z)|} \pi_{k-1}(z) dz.$$

Hence

$$\begin{aligned} \int_{\mathbb{R}^n} \Phi(x) \pi_k(x) dx &= \frac{1}{|D_k|\Phi_{v_k}(y_k)} \int_{\mathbb{R}^n} \frac{\tilde{g}(k, z, y_k)}{|g(k, z)|} \pi_{k-1}(z) dz \\ &= \frac{1}{|D_k|\Phi_{v_k}(y_k)} \\ &\quad \times \int_{\mathbb{R}^n} \int_{\mathbb{R}^n} \frac{\Phi_k(D_k^{-1}(y_k - h(k, x))) \Psi_w(g^{-1}(k, z)(x - f(k, z)))}{|g(k, z)|} \\ &\quad \times \Phi(x) dx \pi_{k-1}(z) dz \\ &= \int_{\mathbb{R}^n} \Phi(x) \left[ \pi_k(x) - \frac{\Phi_k(D_k^{-1}(y_k - h(k, x)))}{|D_k|\Phi_k(y_k)} \right. \\ &\quad \left. \times \int_{\mathbb{R}^n} \pi_{k-1}(z) \frac{\Psi_k(g^{-1}(k, z)(x - f(k, z)))}{|g(k, z)|} dz \right] dx = 0. \end{aligned}$$

Since  $\Phi \in BC(\mathbb{R}^n)$  is arbitrary, the only solution for this equation is

$$\pi_k(x) = \frac{\Phi_k(D_k^{-1}(y_k - h(k, x)))}{|D_k|\Phi_k(y_k)} \int_{\mathbb{R}^n} \frac{\Psi_k(g^{-1}(k, z)(x - f(k, z)))}{|g(k, z)|} \pi_{k-1}(z) dz.$$

At  $k = 0$ ,

$$\begin{aligned} \pi_0(\Phi) &= \bar{E} [\Lambda_0 \Phi(x) | \mathcal{F}_0^y] = \bar{E} \left[ \frac{\Phi_0(D_0^{-1}(y_0 - h(0, x)))}{|D_0|\Phi_0(y_0)} \Phi(x) \right] \\ &= \frac{1}{|D_0|\Phi_0(y_0)} \int_{\mathbb{R}^n} \Phi_0(D_0^{-1}(y_0 - h(0, x))) \Phi(x) d\Pi_{x_0}(x). \end{aligned}$$

Hence,

$$\pi_0(x) = \frac{\Phi_0(D_0^{-1}(y_0 - h(0, x)))}{|D_0|\Phi_0(y_0)} \Pi_{x_0}(x).$$

□

**Example 2.2.20.** (Linear Gaussian filter) It is assumed that the state and observations are given by

$$\left( \Omega, \mathcal{F}, \{\mathcal{G}_k\}, P \right) : \begin{cases} x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1}, & x_k \in \mathbb{R}^n \\ y_k = C_kx_k + D_kv_k, & y_k \in \mathbb{R}^d \end{cases} \quad (2.51)$$

in which the noises are Gaussian distributed as follows,  $w_k \sim N(0, I_n)$ ,  $v_k \sim N(0, I_d)$  while sequences  $\{w_k\}$ ,  $\{v_k\}$  are mutually independent and independent of  $x_0 \sim N(\bar{x}_0, V_0)$ ,  $V_0 > 0$ . The Linearity and Gaussianity implies that

$$\Pi_k(x) = \frac{\pi_k(x)}{\int_{\mathbb{R}^n} \pi_k(x) dx} = \frac{\bar{\alpha}_k(x)}{\int_{\mathbb{R}^n} \bar{\alpha}_k(x) dx}$$

is Gaussian with mean  $\hat{x}_{k|k} = E[x_k | \mathcal{F}_k^y]$  and Variance  $V_{k|k} = E[(x_k - \hat{x}_{k|k})(x_k - \hat{x}_{k|k})^T | \mathcal{F}_k^y]$ .

Recursive equations for  $\hat{x}_{k|k}$ ,  $V_{k|k}$  using the unnormalized conditional density, will be obtained.

According to Theorem 2.2.19,  $\{\bar{\alpha}_k(x)\} \triangleq \{\pi_k(x)\}$  satisfies the following recursion

$$\bar{\alpha}_k(x) = \frac{\Xi_{v_k}(D_k^{-1}(y_k - C_k x))}{|D_k| \Xi_{v_k}(y_k)} \int_{\mathbb{R}^n} \frac{\Psi_{w_k}(B_k^{-1}(x - A_k z))}{|B_k|} \bar{\alpha}_{k-1}(z) dz \quad (2.52)$$

where

$$\Psi_{w_k}(x) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left(-\frac{x^T x}{2}\right), \quad (2.53)$$

$$\Xi_{v_k}(y) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{y^T y}{2}\right). \quad (2.54)$$

A solution to (2.52) having the following form, is assumed

$$\bar{\alpha}_k(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |V_{k|k}|^{\frac{1}{2}}} \exp\left(- (x - \hat{x}_{k|k})^T \frac{(V_{k|k})^{-1}}{2} (x - \hat{x}_{k|k}) + \beta_{k|k}\right) \quad (2.55)$$

when  $\{V_{k|k}, \hat{x}_{k|k}, \beta_{k|k}\}$  will be identified shortly. Substituting (2.55) into the recursion (2.52) deduces the following recursive relations for  $\{V_{k|k}, \hat{x}_{k|k}, \beta_{k|k}\}$ :

$$\hat{x}_{k|k} = V_{k|k} \left[ C_k^T (D_k D_k^T)^{-1} y_k + (V_{k|k-1})^{-1} \hat{x}_{k|k-1} \right], \quad (2.56)$$

$$\hat{x}_{k|k-1} = A_k \hat{x}_{k-1|k-1}, \quad (2.57)$$

$$V_{k|k} = \left( C_k^T (D_k D_k^T)^{-1} C_k + (V_{k|k-1})^{-1} \right)^{-1}, \quad (2.58)$$

$$V_{k|k-1} = B_k B_k^T + A_k V_{k-1|k-1} A_k^T \quad (2.59)$$

$$\begin{aligned} &= (B_k B_k^T) + A_k V_{k-1|k-2} A_k^T \\ &\quad - A_k V_{k-1|k-2} C_{k-1}^T (D_{k-1} D_{k-1}^T + C_{k-1} V_{k-1|k-2} C_{k-1}^T)^{-1} C_{k-1} V_{k-1|k-2} A_k^T, \end{aligned} \quad (2.60)$$



$$\begin{aligned} \beta_{k|k} = & - \sum_{i=1}^k \frac{1}{2} (y_i - C_i \hat{x}_{i|i-1})^T (C_i V_{i|i-1} C_i^T + D_i^T D_i)^{-1} (y_i - C_i \hat{x}_{i|i-1}) + \sum_{i=1}^k \frac{y_i^T y_i}{2} \\ & - \frac{1}{2} \sum_{i=1}^k \log |C_i V_{i|i-1} C_i^T + D_i^T D_i|, \end{aligned} \quad (2.61)$$

with initial conditions

$$\hat{x}_{0|0} = V_{0|0} \left[ C_0^T (D_0 D_0^T)^{-1} y_0 \right], \quad (2.62)$$

$$V_{0|0} = \left( C_0^T (D_0 D_0^T)^{-1} C_0 \right)^{-1}, \quad (2.63)$$

$$\begin{aligned} \beta_{0|0} = & - \frac{1}{2} (y_0 - C_0 \bar{x}_0)^T (C_0 V_0 C_0^T + D_0^T D_0)^{-1} (y_0 - C_0 \bar{x}_0) + \frac{y_0^T y_0}{2} \\ & - \frac{1}{2} \log |C_0 V_0 C_0^T + D_0^T D_0|, \end{aligned} \quad (2.64)$$

thus it is concluded that (2.55) is indeed a solution of the recursion (2.52).

Furthermore,  $\hat{x}_{k|k}$  can be written as

$$\hat{x}_{k|k} = \hat{x}_{k|k-1} + V_{k|k-1} C_k^T \left[ C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right]^{-1} (y_k - C_k \hat{x}_{k|k-1}). \quad (2.65)$$

Also

$$\begin{aligned} \hat{x}_{k|k-1} = & A_k \hat{x}_{k-1|k-2} + A_k V_{k-1|k-2} C_{k-1}^T \left[ C_{k-1} V_{k-1|k-2} C_{k-1}^T + (D_{k-1} D_{k-1}^T) \right]^{-1} \\ & \times (y_{k-1} - C_{k-1} \hat{x}_{k-1|k-2}). \end{aligned} \quad (2.66)$$

This is the solution of the Kalman Filter [1]. See Appendix B for the derivation.



# CHAPTER 3

## ROBUST LEAST-SQUARE ESTIMATION FOR A CLASS OF SYSTEMS

This chapter considers least-square estimation problems for classes of systems. In Section 3.2 the estimation problem for a class of systems is formulated. A general framework is put forward in which the basic ideas and theory are explained and the fundamental results are derived. In Section 3.3 various examples of estimation theory are introduced, when the models (conditional distributions, joint distributions) are uncertain and they belong to specific subsets of the set of conditional or joint distributions, in order to illustrate the results derived. Finally, in Section 3.4 an overview on MIMO communication systems is presented and the theory developed in previous sections is applied to MIMO systems.

### 3.1 Introduction

In classical Least-Square estimation of Random Variables one is interested in finding the best estimate,  $\Phi^*(Y)$ , of a RV  $X$  from the measurements of a RV  $Y$ , by minimizing the expected value of the least-square error  $e(X, \Phi) \triangleq \|X - \Phi(Y)\|_{\mathfrak{R}^n}^2$  over all functions  $\Phi : \mathfrak{R}^d \rightarrow \mathfrak{R}^n, Y \rightarrow \Phi(Y)$ , which is a function of  $Y$  ( $\|x\|_{\mathfrak{R}^n}^2$  denotes Euclidean norm of  $x \in \mathfrak{R}^n$ ). By the orthogonal projection theorem, the solution is given by the conditional expectation:

$$\Phi^*(Y) = E[X|Y] = \int_{\mathfrak{R}^n} x dP_{X|Y}(x|y) \quad (3.1)$$

where  $P_{X|Y}$  is the conditional distribution of  $X$  given  $Y$  [1, 2].

The solution to this estimation problem is stated in terms of the  $\hat{a}$  posteriori distribution function  $P_{X|Y}$ . This distribution contains all information available to estimate  $X$  or any nonlinear function of it. The objective is thus to find, recursively in time, the evaluation equation of the  $\hat{a}$  posteriori distribution and to solve it. The classical result assumes that  $P_{X,Y}$  and hence  $P_{X|Y}$  are completely known. In many real applications of estimation theory, knowledge of  $P_{X,Y}$  and/or  $P_{X|Y}$  is not available; the only knowledge available to the designer is whether  $P_{X,Y}$ , (resp.  $P_{X|Y}$ ) belong to specific classes which are subsets of the set of all joint (resp. conditional) distributions.

This chapter is concerned with estimation techniques, in which the uncertainty description of the system, and the nominal description of the system are modeled by probability distributions, or general measures, defined on measurable spaces. The uncertainty description of these systems is characterized by the class of uncertain measures which satisfy a Kullback-Leibler (KL) distance constraint with respect to a nominal measure. The problem of robust estimation is formulated by minimizing over the set of estimators, the maximum of a linear functional of the uncertain measure over the constraint set. Two type of uncertainty models are considered.

1. Uncertainty Models on Conditional Distributions or otherwise known Stochastic Kernels;
  - i) When the conditional probability distribution of the measurement  $Y$  given the signal to be estimated  $X$ , or channel kernel, is unknown;
  - ii) When the  $\hat{a}$  posteriori distribution of  $X$  given  $Y$  is unknown;
2. Uncertainty Models on Joint Distributions.

Stochastic kernel uncertainty models are appropriate for communication system design, in which the input message has a known distribution, while the channel is unknown but belongs to a certain class of channels. These are nonparametric uncertainty models which so far have not been taken into consideration. Joint distribution uncertainty models are usually employed when both the unobserved and observed random variables are uncertain.

The minimax technique considered here leads to strategies in which the worst case estimate of the uncertain measure subject to the uncertainty description is sought. The theory and contribution of this chapter are developed at two levels of generality; the abstract level and the application level. At the abstract level a general framework

is put forward in which the basic ideas are explained and the fundamental results are derived. Specifically, the estimation problem is described on abstract Polish spaces, while the uncertain model considered is described by stochastic kernels and joint distributions. First, the appropriate space of measures is introduced and then the maximizing kernel and joint measure are computed explicitly using Lagrangian functionals and variational methods. Moreover, important monotonicity properties satisfied by the optimal strategies are presented, which can be used to develop numerical algorithms for computation of the optimal solution, and upper and lower bounds on the optimal solution. At the application level, the results obtained at the abstract level are applied to simple examples.

As was already discussed in Chapter 1 previous related work in which uncertainty is described by relative entropy can be found in [11], [16] and [29]. However, [29] deals only with systems when the uncertainty is defined on joint distributions. On the other hand, [16] similar to [29] deals only with joint distributions, but employs a minimax viewpoint and uses the additional assumption that the system considered is Gaussian. Even though [16] starts with a nonparametric uncertainty models it uses a parametric approach to derive its results. The least-squares error estimator (conditional mean) is considered in [16] and shown to yield a saddle point. A parametric uncertainty model is implemented in [11] which deals with robust least-square error equalization design for multiple-input multiple-output communication channels, when the channel and noise are uncertain. In [11] the Lagrangian duality is used to transform the min-max problem into an equivalent convex min-min problem over a convex domain, to which standard convex optimization methods apply. This way the problem presented in [11] becomes a scalar minimization problem which can be solved numerically. Minimax estimation techniques for uncertain wide sense stationary processes are also considered in [15] and [30], but uncertainty is defined with respect to power spectrum densities. It is worth mentioning that related work in the context of nonlinear stochastic optimal control is also found in [43].

## 3.2 Nonlinear Optimization

This section, formulates and seeks solution to the estimation problem for a class of systems as follows. First, appropriate models for conditional distributions are

introduced. Then appropriate spaces are identified and existence of maximizing measure is shown. Then the theory of Lagrangian functional is invoked to find the maximizing measure. In addition, specific properties of the Lagrangian functional are identified.

### 3.2.1 Formulation on Abstract Spaces

Suppose a measurable space  $(\Omega, \mathcal{F})$  is given on which the unobserved Random Variable (RV),  $X$  and the observed RV  $Y$  are defined, via  $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \Sigma_X)$ ,  $Y : (\Omega, \mathcal{F}) \rightarrow (\mathcal{Y}, \Sigma_Y)$ .

Thus  $\mathcal{X}$  is the space of the unobserved RV, and  $\mathcal{Y}$  is the space of the observed RV. The relation between the unobserved RV  $X$  and the observed RV  $Y$  is defined via a probabilistic mapping. The mapping  $\mu : \mathcal{X} \times \Sigma_Y \rightarrow [0, 1]$  satisfies the following two conditions:

1. For every  $x \in \mathcal{X}$ , the set function  $\mu(x, \cdot)$  is a probability measure on  $\Sigma_Y$  (possibly finite additive);
2. For every  $F \in \Sigma_Y$ , the function  $\mu(\cdot, F)$  is  $\mathcal{X}$ -measurable.

The mapping  $\mu$  is called a stochastic kernel or transition probability. Let  $\mathcal{P}$  denote the class of all stochastic kernels,  $\mathcal{M}_1(\mathcal{X})$  denote the space of probability measures (possibly finite additive) on  $\mathcal{X}$ , and let the measure induced by  $X$ ,  $P_X \in \mathcal{M}_1(\mathcal{X})$  be fixed. For the given pair  $\{\mu, P_X\}$ ,  $\mu \in \mathcal{P}$ ,  $P_X \in \mathcal{M}_1(\mathcal{X})$ , two probability measures may be introduced as follows.

1. The joint probability measure  $P_{X,Y} \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$  defined by

$$P_{X,Y}(G) = (P_X \otimes \mu)(G) \triangleq \int_{\mathcal{X}} \mu(x, G_x) P_X(dx), \quad \forall G \in \Sigma_X \times \Sigma_Y$$

where  $\otimes$  denotes convolution,  $G_x$  is the  $x$ -section of  $G$  defined by  $G_x = \{y \in \mathcal{Y} : (x, y) \in G\}$ .

2. The marginal probability measure  $P_Y \in \mathcal{M}_1(\mathcal{Y})$  corresponding to  $\mu \in \mathcal{P}$  is given by

$$P_Y(F) = P_{X,Y}(\mathcal{X} \times F) \triangleq \int_{\mathcal{X}} \mu(x, F) P_X(dx), \quad \forall F \in \Sigma_Y$$

The objective is to estimate  $X$  by a function of the random variable  $Y$ . Let  $\hat{X} = \Phi(Y)$  denote the estimate of  $X$ . The estimation is done by introducing a pay-off, and then minimizing the pay-off over the class of estimators  $\Phi$  in the admissible set denoted by  $\mathcal{X}_{ad}$ . It is assumed that all admissible estimators  $\Phi : \mathcal{Y} \rightarrow \mathcal{X}$  are  $\Sigma_{\mathcal{Y}}$  measurable and continuous. Let  $\ell : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty)$  be an  $\Sigma_{\mathcal{X}} \times \Sigma_{\mathcal{Y}}$ -measurable function, which corresponds to the sample pay-off. The classical estimation problem deals with minimization of the average pay-off given by

$$J(\Phi^*) = \inf_{\Phi \in \mathcal{X}_{ad}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) dP_{X,Y}(x, y) \quad (3.2)$$

$$= \inf_{\Phi \in \mathcal{X}_{ad}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \mu(x, dy) dP_X(x) \quad (3.3)$$

$$= \inf_{\Phi \in \mathcal{X}_{ad}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \eta(y, dx) dP_Y(y). \quad (3.4)$$

Clearly if  $\ell(x, \Phi(y)) = \|x - \Phi(y)\|_{\mathbb{R}^n}$ , then  $\Phi^*(Y) = E[X|Y]$ . On the other hand, if  $P_{X,Y}$  or  $\mu, \eta$  are unknown, then new estimators should be sought (this is done in Sections 3.2.2, 3.2.3 and 3.2.4). It is noted that (3.2) will be used when the uncertainty is on the joint distribution, while (3.3), (3.4) will be used when the uncertainty is on the channel kernel, *a posteriori* distribution, respectively.

Next the appropriate topologies and function spaces used in this paper are introduced. Throughout the rest of this chapter it is assumed that both  $\mathcal{X}$  and  $\mathcal{Y}$  are Polish spaces (complete separable metric spaces) and therefore normal topological spaces [45]. Let  $BC(\mathcal{Y})$  denote the vector space of bounded continuous real valued functions defined on the Polish space  $\mathcal{Y}$ . Furnished with the sup norm topology, this is a Banach space. Let  $(BC(\mathcal{Y}))^*$  denote its topological dual. It is known [45, pg. 262] that  $(BC(\mathcal{Y}))^*$  is isometrically isomorphic to the Banach space of finitely additive regular bounded signed measures on  $\Sigma_{\mathcal{Y}}$ . Denote this by  $M_{rba}(\mathcal{Y})$  and let  $\Pi_{rba}(\mathcal{Y}) \subset M_{rba}(\mathcal{Y})$  denote the set of regular bounded finitely additive probability measures on  $\mathcal{Y}$ . Clearly if  $\mathcal{Y}$  is compact, then  $(BC(\mathcal{Y}))^*$  will be the space of countably additive signed measures. Let  $L_1(P_X, BC(\mathcal{Y}))$  denote the space of all  $P_X$  integrable functions defined on  $\mathcal{X}$  with values in  $BC(\mathcal{Y})$ . In other words, for each  $\phi \in L_1(P_X, BC(\mathcal{Y}))$

$$\|\phi\|_{P_X} \equiv \int_{\mathcal{X}} \|\phi(x)(\cdot)\|_{BC(\mathcal{Y})} P_X(dx) < \infty.$$

With respect to this norm topology this is a Banach space. Since the Banach spaces  $BC(\mathcal{Y})$  and its dual  $M_{rba}(\mathcal{Y})$  do not satisfy the Radon Nikodym property, the dual

of  $L_1(P_X, BC(\mathcal{Y}))$  is not  $L_\infty(P_X, M_{rba}(\mathcal{Y}))$ . However, it follows from the theory of "lifting" [46, Theorem 7, pg. 94; Theorem 9, pg. 97] that the dual of  $L_1(P_X, BC(\mathcal{Y}))$  is  $L_\infty^w(P_X, M_{rba}(\mathcal{Y}))$ , i.e., the space of all  $M_{rba}(\mathcal{Y})$  valued functions  $\{\mu\}$  which are weak star measurable, in the sense that for each  $\phi \in BC(\mathcal{Y})$ ,  $x \rightarrow \mu_x(\phi) \equiv \int_{\mathcal{Y}} \phi(z) \mu(x, dz)$  is  $P_X$  measurable and  $P_X$ -essentially bounded. Now define the admissible set as follows

$$\mathcal{P} \equiv L_\infty^w(P_X, \Pi_{rba}(\mathcal{Y})) \subset L_\infty^w(P_X, M_{rba}(\mathcal{Y})).$$

In other words,  $\mathcal{P}$  is the unit sphere in the space  $L_\infty^w(P_X, M_{rba}(\mathcal{Y}))$ . It is assumed throughout this chapter that  $\ell : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty]$  is  $\Sigma_{\mathcal{X}} \times \Sigma_{\mathcal{Y}}$  measurable (sample pay-off) function from the class  $L_1(P_X, BC(\mathcal{Y}))$ .

### 3.2.2 Uncertainty on the Channel Kernel and Minimax Pay-off

It is assumed that the probabilistic kernel  $\mu \in \mathcal{P}$  introduced earlier represents the nominal system model or mapping, which is fixed. The true kernel denoted by  $\nu : \mathcal{X} \times \Sigma_{\mathcal{Y}} \rightarrow [0, 1]$  is assumed unknown. Envisioned scenarios are communication channels whose nominal behavior is known, while its true conditional distribution is unknown. The KL distance will be used as a measure of distance between the true model and uncertainty model, hence the true kernel is assumed to belong to the pointwise uncertainty set defined by

$$\mathcal{A}^x(\mu) \triangleq \left\{ \nu \in \mathcal{P} : H(\nu|\mu)(x) \leq R(x) \right\} \quad (3.5)$$

where  $R : \mathcal{X} \rightarrow [0, \infty)$  and  $H(\cdot|\cdot)(x) : \mathcal{X} \rightarrow [0, \infty]$  is the KL distance between two kernels defined by

$$H(\nu|\mu)(x) = \begin{cases} \int_{\mathcal{Y}} \log \frac{\nu(x, dy)}{\mu(x, dy)} \nu(x, dy), & \text{if } \nu(x, \cdot) \ll \mu(x, \cdot), P_X - a.s. \\ \infty, & \text{otherwise.} \end{cases}$$

It is also assumed that  $R \in BC(\mathcal{X})$ . Moreover,  $\mu : \mathcal{X} \times \Sigma_{\mathcal{Y}} \rightarrow [0, 1]$  is the nominal fixed kernel.

**Remark 3.2.1.** *The value of  $R(x)$  can be evaluated by statistical methods as follows. Using experimental data and counting techniques the different possible conditional distributions  $\nu(x, dy)$  can be found. On the other hand, if the true distribution  $\nu(x, dy)$  is parameterized then by using counting techniques in finding relative entropy, the function  $R(x)$  is determined.*



Additionally, the following uncertainty set is defined

$$\mathcal{A}(\mu) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{X}} H(\nu|\mu)(x) dP_X(x) \leq \int_{\mathcal{X}} R(x) dP_X(x) \stackrel{\nabla}{=} r_1 \right\}. \quad (3.6)$$

It is assumed that  $\mathcal{A}(\mu)$  is non-empty. Since, the true kernel  $\nu \in \mathcal{A}(\mu)$ , is unknown, then the estimation problem can be formulated as a minimax problem defined by

$$\begin{aligned} J_1(\Phi, \nu) &\triangleq \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x), \\ J_1(\Phi^*, \nu^*) &= \inf_{\Phi \in \mathcal{X}_{ad}} \sup_{\nu \in \mathcal{A}(\mu)} J_1(\Phi, \nu). \end{aligned} \quad (3.7)$$

Next the issue of the existence of a solution to the above problem is discussed. Let  $\ell : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty]$  be any  $\Sigma_{\mathcal{X}} \times \Sigma_{\mathcal{Y}}$  measurable function from the class  $L_1(P_X, BC(\mathcal{Y}))$ .

The set  $\mathcal{P}$  is  $w^*$ -compact. This follows from the Alaoglu's theorem [45, Theorem V.4.2, pg. 424]. Also using the lower semi continuity property of relative entropy [47, Lemma 1.4.3, pg. 36]<sup>1</sup>, it follows that the set  $\mathcal{A}(\mu)$  is a  $w^*$ -compact (as a  $w^*$ -closed subset of the  $w^*$ -compact set  $\mathcal{P}$ ). The next lemma establishes the upper semi-continuity of the cost function.

**Lemma 3.2.2.** *Let  $\mathcal{X}, \mathcal{Y}$  be two Polish spaces and  $\ell : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty]$ , a measurable, nonnegative, extended real valued function, and  $y \rightarrow \ell(x, y)$  be continuous, for  $P_X$ -almost all  $x \in \mathcal{X}$ . Also assume  $\Phi : \mathcal{Y} \rightarrow \mathcal{X}$  is  $\Sigma_{\mathcal{Y}}$  measurable and continuous. Then the mapping  $\nu \rightarrow \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) P_X(dx)$  is upper semi-continuous in the  $w^*$ -sense.*

*Proof.* Let  $\{\nu_\alpha\} \in \mathcal{A}(\mu)$  be a net, where  $\alpha \in (\mathcal{D}, \leq)$  (a directed set). Since  $\mathcal{A}(\mu)$  is weak star compact, there exists a subnet of the net  $\{\nu_\alpha\}$ , relabeled as the original net, and an element  $\nu \in \mathcal{A}(\mu)$  such that  $\nu_\alpha \xrightarrow{w^*} \nu$ .<sup>2</sup>

$$\begin{aligned} \overline{\lim}_\alpha \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \wedge m \nu_\alpha(x, dy) P_X(dx) &= \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \wedge m \nu(x, dy) P_X(dx) \\ &\leq \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) P_X(dx) \end{aligned}$$

<sup>1</sup>Examining the proof in [47] one can easily verify that the same procedure holds true not only for countably additive measures but also for finitely additive ones.

<sup>2</sup>That is,

$$\lim_\alpha \left| \int_{\mathcal{X}} \int_{\mathcal{Y}} \phi(x, y) \nu_\alpha(x, dy) P_X(dx) - \int_{\mathcal{X}} \int_{\mathcal{Y}} \phi(x, y) \nu(x, dy) P_X(dx) \right| = 0$$

for any  $\phi \in L_1(P_X; BC(\mathcal{Y}))$ .

where  $m \in N$  ( $N$  denotes non-negative integers) is arbitrary. Then there exists  $\alpha_0 \in \mathcal{D}$  such that

$$\sup_{\alpha \geq \alpha_0} \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \wedge m \nu_{\alpha}(x, dy) P_X(dx) \leq \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) P_X(dx). \quad (3.8)$$

On the other hand, the following holds

$$\begin{aligned} \overline{\lim}_{\alpha} \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \nu_{\alpha}(x, dy) P_X(dx) &= \overline{\lim}_{\alpha} \sup_{m \in N} \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \wedge m \nu_{\alpha}(x, dy) P_X(dx) \\ &\leq \sup_{\alpha \geq \alpha_0} \sup_{m \in N} \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \wedge m \nu_{\alpha}(x, dy) P_X(dx) \\ &\leq \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) P_X(dx) \end{aligned}$$

where the last inequality follows from (3.8) and the fact that two supremums can be interchanged.

Thus, since  $\mathcal{A}(\mu)$  is  $w^*$ -compact, and

$$\nu \rightarrow \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) P_X(dx)$$

is upper semicontinuous, from the Weierstrass theorem, existence of solution to the supremum problem in (3.7) is established.  $\square$

**Remark 3.2.3.** Note that by a generalization of von Neumann's minimax theorem, the problem (3.7) with  $\mathcal{X}_{ad}$  convex,  $\ell(x, \Phi) = \|x - \Phi\|_{\mathbb{R}^n}$ , satisfies (see Theorem 2.1.15)

$$\begin{aligned} &\inf_{\Phi \in \mathcal{X}_{ad}} \max_{\nu \in \mathcal{A}(\mu)} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x) \\ &= \max_{\nu \in \mathcal{A}(\mu)} \min_{\Phi \in \mathcal{X}_{ad}} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x) \end{aligned} \quad (3.9)$$

since

- 1)  $\mathcal{A}(\mu)$  is compact,
- 2)  $J_1(\Phi, \nu)$  as a function on  $\mathcal{X}_{ad} \times \mathcal{A}(\mu)$  is convex-concave, and
- 3)  $J_1(\Phi, \nu)$  as a function on  $\mathcal{X}_{ad} \times \mathcal{A}(\mu)$  is upper-semicontinuous in  $\nu$  for each  $\Phi \in \mathcal{X}_{ad}$ .

Clearly, for the case of this chapter  $\mathcal{A}(\mu)$  is convex (because of the convexity of relative entropy [47]),  $\mathcal{A}(\mu)$  is  $weak^*$  compact (see statement above Lemma 3.2.2),  $J_1(\Phi, \nu)$  is convex in  $\mathcal{X}_{ad}$  and concave in  $\mathcal{A}(\mu)$ , and  $J_1(\Phi, \nu)$  is upper-semicontinuous in  $\nu$  for each  $\Phi \in \mathcal{X}_{ad}$  (by Lemma 3.2.2). Therefore, the only additional assumption needed for existence of saddle point (hence  $\inf$  becomes  $\min$ ) is that A1)  $\mathcal{X}_{ad}$  is compact and convex, A2)  $J_1(\cdot, \nu)$  is continuous in  $\Phi$  for each  $\nu \in \mathcal{A}(\mu)$  and hence lower-semicontinuous.

Define

$$J_1(\phi, \nu^*) \triangleq \sup_{\nu \in \mathcal{A}(\mu)} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \phi(y)) \nu(x, dy) dP_X(x).$$

The solution of the supremum over  $\mathcal{A}(\mu)$  is resolved by introducing a pair of Lagrange multipliers  $\{\lambda(x), s\}$ ,  $s \in \mathfrak{R}$ ,  $\lambda \in L_1(P_X)$  and defining the Lagrangian functional

$$\begin{aligned} L_1(\nu, \lambda, s) \triangleq & \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x) - s \int_{\mathcal{X}} (H(\nu|\mu)(x) - R(x)) dP_X(x) \\ & - \int_{\mathcal{X}} \lambda(x) \left( \int_{\mathcal{Y}} \nu(x, dy) - 1 \right) dP_X(x) \end{aligned} \quad (3.10)$$

and the dual functional

$$L_1(\nu^*, \lambda^*, s^*) = \inf_{s \geq 0} \inf_{\{\lambda \in L_1(P_X), \lambda \geq 0\}} \sup_{\nu \in \mathcal{P}} L_1(\nu, \lambda, s). \quad (3.11)$$

Alternatively, if the uncertainty set  $\mathcal{A}^x(\mu)$  is considered, the supremum problem is defined as

$$\tilde{J}_1(\Phi, \nu^*) = \sup_{\nu \in \mathcal{A}^x(\mu)} \int_{\mathcal{X}} \int_{\mathcal{Y}} \ell(x, \phi(y)) \nu(x, dy) dP_X(x).$$

In this case the solution of the supremum over  $\mathcal{A}^x(\mu)$  is resolved using lagrange multiplier  $\{\lambda(x), \tilde{s}\}$ ,  $\lambda \in L_1(P_X)$ ,  $\tilde{s} \in L_1(P_X)$ , and defining the Lagrangian functional

$$\begin{aligned} \tilde{L}_1(\nu, \lambda, \tilde{s}) \triangleq & \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x) - \int_{\mathcal{X}} \tilde{s}(x) (H(\nu|\mu)(x) - R(x)) dP_X(x) \\ & - \int_{\mathcal{X}} \lambda(x) \left( \int_{\mathcal{Y}} \nu(x, dy) - 1 \right) dP_X(x) \end{aligned} \quad (3.12)$$

The dual functional is given by

$$\tilde{L}_1(\nu^*, \lambda^*, \tilde{s}^*) = \inf_{\{\tilde{s} \in L_1(P_X), \tilde{s} \geq 0\}} \inf_{\{\lambda \in L_1(P_X), \lambda \geq 0\}} \sup_{\nu \in \mathcal{P}} \tilde{L}_1(\nu, \lambda, \tilde{s}) \quad (3.13)$$

Now, if the infimum over the functions  $\tilde{s} \in L_1(P_X)$  is taken over the real numbers, then the following inequality holds

$$\tilde{L}_1(\nu^*, \lambda^*, \tilde{s}^*) \leq \inf_{\tilde{s} \geq 0} \inf_{\{\lambda \in L_1(P_X), \lambda \geq 0\}} \sup_{\nu \in \mathcal{P}} \tilde{L}_1(\nu, \lambda, \tilde{s}) = L_1(\nu^*, \lambda, \tilde{s}) \quad (3.14)$$

Hence, uncertainty modeling using  $\mathcal{A}(\mu)$  yields higher pay-off than uncertainty modeling using  $\mathcal{A}^x(\mu)$ .

Next, the equivalence between the constrained problem  $J_1(\Phi, \nu^*)$  and the unconstrained problem  $L_1(\nu^*, \lambda^*, s^*)$  is established.

**Theorem 3.2.4.** Suppose  $\ell : \mathcal{X} \times \mathcal{X} \longrightarrow \overline{\mathbb{R}}_0 \equiv [0, \infty]$  is continuous in the second argument, and  $\Phi : \mathcal{Y} \rightarrow \mathcal{X}$  is continuous. If for a given  $\Phi \in \mathcal{X}_{ad}$ ,  $\sup_{\nu \in \mathcal{A}(\mu)} J_1(\Phi, \nu)$  is finite, then the constrained problem (inner supremum) in (3.7), is equivalent to an unconstrained problem as stated below:

$$\sup_{\nu \in \mathcal{A}(\mu)} J_1(\Phi, \nu) = \inf_{s \geq 0} \inf_{\lambda(x)} \sup_{\nu \in \mathcal{A}(\mu)} L_1(\nu, \lambda, s).$$

Further the infimum occurs on the boundary of the set  $\mathcal{A}(\mu)$ , that is

$$\int_{\mathcal{X}} H(\nu^* | \mu)(x) dP_X(x) \Big|_{s=s^*} = r_1.$$

*Proof.* The proof is based on Lagrange Duality theorem [40, Theorem 1, pg. 224].  $\overline{\mathcal{X}}$  is chosen as  $\overline{\mathcal{X}} \equiv L_{\infty}^w(P_X, M_{rba}(\mathcal{Y}))$  which is clearly a vector space. For the set  $\Omega$  the natural choice is the set  $\Omega = L_{\infty}^w(P_X, \Pi_{rba}(\mathcal{Y})) \subseteq \overline{\mathcal{X}}$ . Clearly,  $\int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x)$  is a linear functional of  $\nu$ . Theorem 1, in [40, pg. 224] deals with minimization of a convex functional. Multiplying the equation (4) [40, pg. 224] by a minus sign, converts the problem to maximization of a concave functional over the set  $\Omega$ . This can be applied to maximization of  $J_1(\Phi, \nu)$  over the constrained set, since  $\int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(x, dy) dP_X(x)$  is a linear functional of  $\nu$  (hence concave). Also, it follows that  $\Omega$  is a convex set. Take  $G(\nu) = \int_{\mathcal{X}} H(\nu | \mu)(x) - R(x) dP_X(x)$ , where  $\nu \in \overline{\mathcal{X}}$ . Then  $G : \overline{\mathcal{X}} \rightarrow \mathfrak{R}$ , is a convex mapping from  $\overline{\mathcal{X}}$  into the ordered vector space  $(\mathfrak{R}, \prec)$  with natural ordering. If  $\nu = \mu$  is chosen, then  $G(\nu) = - \int_{\mathcal{X}} R(x) dP_X(x) < 0$ <sup>3</sup>. Hence there exists a measure  $\nu \in \overline{\mathcal{X}}$  such that  $G(\nu) < 0$ . So when  $\sup_{\nu \in \mathcal{A}(\mu)} J_1(\Phi, \nu)$  is finite, conditions of the theorem are satisfied, and the constrained and unconstrained problems are equivalent. Also according to the same duality theorem, if the supremum is achieved by some  $\nu^* \in L_{\infty}^w(P_X, M_{rba}(\mathcal{Y}))$ , then

$$s \left( \int_{\mathcal{X}} (H(\nu^* | \mu)(x) - R(x)) dP_X(x) \right) = 0.$$

In other words, for non-zero  $s \in (0, \infty)$ , the solution occurs on the boundary.  $\square$

The solution of the maximization over  $\nu \in \mathcal{A}(\mu)$  of (3.10) is presented in the next Theorem.

<sup>3</sup>The inequality is strict if  $R$  is non-zero on a set of non-zero measure, i.e., if there exists  $E \in \Sigma_{\mathcal{X}}$ , such that  $P_X(E) \neq 0$ , and  $R(x) \neq 0$  for  $x \in E$ . If such a set does not exist, then  $R$  would be zero  $P_X$  almost everywhere, and the problem becomes trivial.

**Theorem 3.2.5.** *Suppose the condition of Theorem 2.1.16 holds. The supremum of (3.11) over  $\nu \in \mathcal{A}(\mu)$  is given by*

$$\nu^*(x, dy) = \frac{e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy)}{\int_{\mathcal{Y}} e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy)} \quad (3.15)$$

where  $s \geq 0$  is found by the constraint

$$\int_{\mathcal{X}} H(\nu^* | \mu)(x) dP_X(x) \Big|_{s=s^*} = r_1. \quad (3.16)$$

Moreover,

$$L_1(\nu^*, \lambda^*, s^*) = \inf_{s \geq 0} \int_{\mathcal{X}} s \log \int_{\mathcal{Y}} e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy) dP_X(x) + s \int_{\mathcal{X}} R(x) dP_X(x). \quad (3.17)$$

*Proof.* The Gateaux derivative of  $L_1(\cdot, \lambda, s)$  at any  $\nu^*$  in the direction  $\nu - \nu^*$  is derived by computing the following expression

$$\frac{d}{d\varepsilon} L_1(\nu^* + \varepsilon(\nu - \nu^*), \lambda, s) \Big|_{\varepsilon=0} \triangleq \delta L_1(\nu^*; \nu - \nu^*). \quad (3.18)$$

After some calculations it is deduced to

$$\begin{aligned} \delta L_1(\nu^*; \nu - \nu^*) &= \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) (\nu(x, dy) - \nu^*(x, dy)) dP_X(x) \\ &\quad - \int_{\mathcal{X}} s \left( \int_{\mathcal{Y}} \log \frac{\nu^*(x, dy)}{\mu(x, dy)} (\nu(x, dy) - \nu^*(x, dy)) \right) dP_X(x) \\ &\quad - \int_{\mathcal{X}} \lambda(x) \int_{\mathcal{Y}} (\nu(x, dy) - \nu^*(x, dy)) dP_X(x) \\ &= \int_{\mathcal{X} \times \mathcal{Y}} \log \left( e^{\ell(x, \Phi(y)) - \lambda(x)} \left( \frac{\nu^*(x, dy)}{\mu(x, dy)} \right)^{-s} \right) (\nu(x, dy) \\ &\quad - \nu^*(x, dy)) dP_X(x). \end{aligned} \quad (3.19)$$

Since  $L_1(\nu, \lambda, s)$  is concave in  $\nu$ , then from basic principles of calculus of variation a necessary and sufficient condition for  $\nu^*$  to be the maximizer measure is that  $\delta L_1(\nu^*; \nu - \nu^*) = 0, \forall \nu \in \mathcal{P}$ . Since the inequality holds for all  $\nu$  the Gateaux gradient must vanish, hence

$$\left( \frac{\nu^*(x, dy)}{\mu(x, dy)} \right)^{-s} = e^{-\ell(x, \Phi(y)) + \lambda(x)}.$$

Moreover,

$$\frac{\nu^*(x, dy)}{\mu(x, dy)} = e^{\frac{\ell(x, \Phi(y)) - \lambda(x)}{s}}.$$

Since  $\nu^*(x, dy)$  is a probability measure on  $\mathcal{Y}$ , then using the constraint  $\int_{\mathcal{Y}} \nu^*(x, dy) = 1$ , the following function is obtained

$$e^{-\frac{\lambda(x)}{s}} = \frac{1}{\int_{\mathcal{Y}} e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy)}.$$

Thus,

$$\nu^*(x, dy) = \frac{e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy)}{\int_{\mathcal{Y}} e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy)}, \quad \forall x \in \mathcal{X}. \quad (3.20)$$

Substituting  $\nu^*$  into (3.10) yields

$$L_1(\nu^*, \lambda, s) = \int_{\mathcal{X}} s \log \int_{\mathcal{Y}} e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy) dP_X(x) + \int_{\mathcal{X}} s R(x) dP_X(x). \quad (3.21)$$

□

### 3.2.3 Uncertainty on the $\hat{a}$ Posteriori Distribution and Minimax Pay-off

In this section an uncertainty model on the  $\hat{a}$  posteriori distribution is considered. For this case, one may consider the mapping  $\eta : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$  that satisfies the following two conditions:

1. For every  $y \in \mathcal{Y}$ , the set function  $\eta(y, \cdot)$  is a probability measure on  $\Sigma_{\mathcal{X}}$  (possibly finite additive);
2. For every  $F \in \Sigma_{\mathcal{X}}$ , the function  $\eta(\cdot, F)$  is  $\mathcal{Y}$ -measurable.

This probabilistic kernel  $\eta(y, dx)$  represents the nominal system model or mapping ( $\hat{a}$  posteriori information). The true kernel  $\nu(y, dx)$ , denoted by  $\nu : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$  belongs to an uncertainty set described by

$$\mathcal{B}^y(\eta) \triangleq \left\{ \nu \in \mathcal{P} : H(\nu|\eta)(y) \leq R(y) \right\}$$

where  $R : \mathcal{Y} \rightarrow [0, \infty)$  and  $H(\cdot|\cdot)(y) : \mathcal{Y} \rightarrow [0, \infty]$  is the KL distance between two kernels defined, in a similar way as before, by

$$H(\nu|\eta)(x) = \begin{cases} \int_{\mathcal{X}} \log \frac{\nu(y, dx)}{\eta(y, dx)} \nu(y, dx), & \text{if } \nu(y, \cdot) \ll \eta(y, \cdot), P_Y - a.s. \\ \infty, & \text{otherwise.} \end{cases}$$

It is assumed that  $R \in BC(\mathcal{Y})$ . Moreover,  $\eta : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$  is the nominal fixed kernel. The value of  $R(y)$  can be evaluated by statistical methods as explained in Remark 3.2.1.

It is also assumed that  $\mathcal{B}^y(\eta)$  is non-empty. Since, the true kernel  $\nu \in \mathcal{B}^y(\eta)$ , is unknown, then the estimation problem can be formulated as a minimax problem defined by

$$J_2(\Phi^*, \nu^*) = \inf_{\Phi \in \mathcal{X}_{ad}} \sup_{\nu \in \mathcal{B}^y(\eta)} \int_{\mathcal{X}} \ell(x, \Phi(y)) \nu(y, dx). \quad (3.22)$$

Existence and equivalence of constrained and unconstrained problem is shown similarly using Theorem 2.1.16 following the same procedure as in Section 3.2.2.

The solution of the supremum over  $\mathcal{B}^y(\eta)$  is resolved by introducing a pair of Lagrange multipliers  $\{\lambda(y), \tilde{s}(y)\}$ ,  $\lambda, \tilde{s} \in L_1(P_Y)$  and defining the Lagrangian

$$\begin{aligned} L_2(\nu, \lambda, \tilde{s}) \triangleq & \int_{\mathcal{X}} \ell(x, \Phi(y)) \nu(y, dx) - \tilde{s}(y) \left( H(\nu|\eta)(y) - R(y) \right) \\ & - \lambda(y) \left( \int_{\mathcal{X}} \nu(y, dx) - 1 \right) \end{aligned} \quad (3.23)$$

and the dual functional

$$L_2(\nu^*, \lambda^*, \tilde{s}^*) = \inf_{\{\tilde{s} \in L_1(P_Y), \tilde{s} \geq 0\}} \inf_{\{\lambda \in L_1(P_Y), \lambda \geq 0\}} \sup_{\nu \in \mathcal{P}} L_2(\nu, \lambda, \tilde{s}). \quad (3.24)$$

The solution of the maximization over  $\mathcal{B}^y(\eta)$  of (3.23) is presented in the next Theorem.

**Theorem 3.2.6.** *Suppose the condition of Theorem 2.1.16 holds. The supremum of (3.24) over  $\mathcal{B}^y(\eta)$  is given by*

$$\nu^*(y, dx) = \frac{e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \eta(y, dx)} \quad (3.25)$$

where  $\tilde{s} \in L_1(P_Y)$  is found by the constraint

$$H(\nu^*|\eta)(y) \Big|_{\tilde{s}(y)=\tilde{s}^*(y)} = R(y). \quad (3.26)$$

Moreover,

$$L_2(\nu^*, \lambda^*, \tilde{s}^*) = \inf_{\{\tilde{s} \in L_1(P_Y), \tilde{s} \geq 0\}} \tilde{s}(y) \log \left( \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \eta(y, dx) \right) + \tilde{s}(y) R(y). \quad (3.27)$$

*Proof.* The proof is similar with the one in section 3.2.2. The Gateaux derivative of  $L_2(\cdot, \lambda, \tilde{s})$  at any  $\nu^*$  in the direction  $\nu - \nu^*$  is derived by computing the following expression

$$\left. \frac{d}{d\varepsilon} L_1(\nu^* + \varepsilon(\nu - \nu^*), \lambda, \tilde{s}) \right|_{\varepsilon=0} \triangleq \delta L_1(\nu^*; \nu - \nu^*). \quad (3.28)$$

After some calculations the following expression is deduced

$$\begin{aligned} \delta L_2(\nu^*; \nu - \nu^*) &= \int_{\mathcal{X}} \ell(x, \Phi(y)) (\nu(y, dx) - \nu^*(y, dx)) \\ &\quad - \tilde{s}(y) \left( \int_{\mathcal{X}} \log \frac{\nu^*(y, dx)}{\eta(y, dx)} (\nu(y, dx) - \nu^*(y, dx)) \right) \\ &\quad - \lambda(y) \int_{\mathcal{X}} (\nu(y, dx) - \nu^*(y, dx)) \\ &= \int_{\mathcal{X}} \log \left( e^{\ell(x, \Phi(y)) - \lambda(y)} \left( \frac{\nu^*(y, dx)}{\eta(y, dx)} \right)^{-\tilde{s}(y)} \right) (\nu(y, dx) - \nu^*(y, dx)). \end{aligned} \quad (3.29)$$

Since  $L_2(\nu, \lambda, \tilde{s})$  is concave in  $\nu$ , then from basic principles of calculus of variation a necessary and sufficient condition for  $\nu^*$  to be the maximizer measure is that  $\delta L_2(\nu^*; \nu - \nu^*) = 0, \forall \nu \in \mathcal{P}$ . Since the inequality holds for all  $\nu$  the Gateaux gradient must vanish, hence

$$\left( \frac{\nu^*(x, dy)}{\mu(x, dy)} \right)^{-\tilde{s}(y)} = e^{-\ell(x, \Phi(y)) + \lambda(y)}.$$

Moreover,

$$\frac{\nu^*(x, dy)}{\eta(x, dy)} = e^{\frac{\ell(x, \Phi(y)) - \lambda(y)}{\tilde{s}(y)}}.$$

Since  $\nu^*(y, dx)$  is a probability measure on  $\mathcal{X}$ , then using the constraint  $\int_{\mathcal{X}} \nu^*(y, dx) = 1$ , the following function is obtained

$$e^{-\frac{\lambda(y)}{\tilde{s}(y)}} = \frac{1}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \eta(y, dx)}.$$

Thus,

$$\nu^*(y, dx) = \frac{e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \eta(y, dx)}, \quad \forall y \in \mathcal{Y}. \quad (3.30)$$

Substituting  $\nu^*$  into (3.23) yields

$$L_2(\nu^*, \lambda, \tilde{s}) = \inf_{\{\tilde{s} \in L_1(P_Y), \tilde{s} \geq 0\}} \tilde{s}(y) \log \left( \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \eta(y, dx) \right) + \tilde{s}(y) R(y). \quad (3.31)$$

□



Alternatively, one may consider the constraint

$$\mathcal{B}(\eta) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{Y}} H(\nu|\eta)(y) dP_Y(y) \leq \int_{\mathcal{Y}} R(y) dP_Y(y) \stackrel{\nabla}{=} r_2 \right\}.$$

This estimation problem can be formulated as a minimax problem defined by

$$J_3(\Phi^*, \nu^*) = \inf_{\Phi \in \mathcal{X}_{ad}} \sup_{\nu \in \mathcal{B}(\eta)} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(y, dx) dP_Y(y). \quad (3.32)$$

Defining the Lagrangian as before (with the Lagrange multiplier  $\tilde{s}$  replaced by a real number); then its supremum over  $\nu \in \mathcal{P}$  is given by

$$\nu^*(y, dx) = \frac{e^{\frac{\ell(x, \Phi(y))}{s}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s}} \eta(y, dx)} \quad (3.33)$$

where  $s \geq 0$  is found by the constraint

$$\int_{\mathcal{Y}} H(\nu^*|\eta)(y) dP_Y(y) \Big|_{s=s^*} = r_2. \quad (3.34)$$

Moreover,

$$L_3(\nu^*, \lambda^*, s^*) = \inf_{s \geq 0} \int_{\mathcal{Y}} s \log \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s}} \eta(y, dx) dP_Y(y) + \int_{\mathcal{Y}} s R(y) dP_Y(y). \quad (3.35)$$

**Remark 3.2.7.** Next, a new case is investigated, when the true kernel  $\nu(y, dx)$ , denoted by  $\nu : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$  belongs to a new uncertainty set described by

$$\mathcal{B}^R(\eta) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{Y}} H(\nu|\eta)(y) dP_Y(y) \leq \int_{\mathcal{X} \times \mathcal{Y}} R(x, y) \nu(y, dx) dP_Y(y) + \bar{R} \stackrel{\nabla}{=} r_3 \right\},$$

where  $R : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \infty)$ ,  $R \in BC(\mathcal{X} \times \mathcal{Y})$ . The estimation problem can be formulated using the uncertainty set  $\mathcal{B}^*(\eta)$  as follows.

$$J_3^R(\Phi^*, \nu^*) = \inf_{\Phi \in \mathcal{X}_{ad}} \sup_{\nu \in \mathcal{B}^R(\eta)} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) \nu(y, dx) dP_Y(y). \quad (3.36)$$

Defining the Lagrangian as before, then its supremum over  $\nu \in \mathcal{P}$  is given by

$$\nu^*(y, dx) = \frac{e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \eta(y, dx)} \quad (3.37)$$

where  $s \geq 0$  is found by the constraint

$$\int_{\mathcal{Y}} H(\nu^*|\eta)(y) dP_Y(y) \Big|_{s=s^*} = r_3. \quad (3.38)$$

Moreover,

$$L_3^R(\nu^*, \lambda^*, s^*) = \inf_{s \geq 0} \int_{\mathcal{Y}} s \log \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \eta(y, dx) dP_Y(y) + s \bar{R}. \quad (3.39)$$

### 3.2.4 Uncertainty on the Joint Distribution and Minimax Pay-off

In this section, the uncertainty is modeled via the joint distribution of  $X$  and  $Y$ . This model is appropriate when one wishes to model *a priori* uncertainty. The pay-off function (3.2) is considered

$$J_4(\Phi, Q_{X,Y}) = \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) dQ_{X,Y}(x, y)$$

in which  $Q_{X,Y}$  is the uncertain joint measure. Assume the joint distribution  $P_{X,Y}(x, y)$  represents the nominal system model. The true joint distribution denoted by  $Q \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$ , is assumed to belong to an uncertainty set described by

$$\mathcal{C}(P_{X,Y}) \triangleq \{Q_{X,Y} \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y}) : H(Q_{X,Y}|P_{X,Y}) \leq R\}$$

where  $R \in [0, \infty)$  and  $H(\cdot|\cdot)$  is the KL distance between the two joint distributions defined by

$$H(Q_{X,Y}|P_{X,Y}) = \begin{cases} \int_{\mathcal{X} \times \mathcal{Y}} \log \frac{Q_{X,Y}(x,y)}{P_{X,Y}(x,y)} Q_{X,Y}(x, dy) & \text{if } Q_{X,Y} \ll P_{X,Y} \\ \infty, & \text{otherwise.} \end{cases}$$

The value of  $R$  can be evaluated by statistical methods as explained in Remark 3.2.1.

Since the true probability measure  $Q_{X,Y} \in \mathcal{C}(P_{X,Y})$ , then the estimation problem can be formulated as a minimax problem defined by

$$J_4(\Phi^*, Q_{X,Y}^*) = \inf_{\Phi \in \mathcal{X}_{ad}} \sup_{Q_{X,Y} \in \mathcal{C}(P_{X,Y})} \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) dQ_{X,Y}(x, y). \quad (3.40)$$

This type of problem is similar to the one pursued in [16] and [29], although the method and results presented here are different. The system considered here is an abstract nonparametric system which can be used for various problems, where as in [16] a specific gaussian system is considered and the results presented are base on a parametric uncertainty model. In [29], even though a nonparametric model is used, the method described does not include a minimax approach and it is only shown that risk-sensitive estimators enjoy an error bound.

Existence and equivalence of constrained and unconstrained problem is shown similarly using Theorem 2.1.16. The supremum in (3.40) is resolved by defining the Lagrangian

$$L_4(Q_{X,Y}^*, \lambda^*, s^*) \triangleq \inf_{s \geq 0} \inf_{\lambda \geq 0} \sup_{Q_{X,Y} \in \mathcal{C}(P_{X,Y})} \left\{ \int_{\mathcal{X} \times \mathcal{Y}} \ell(x, \Phi(y)) dQ_{X,Y}(x, y) \right.$$

$$-s \left( H(Q_{X,Y}|P_{X,Y}) - R \right) - \lambda \left( \int_{\mathcal{X} \times \mathcal{Y}} dQ_{X,Y}(x, y) - 1 \right) \}. \quad (3.41)$$

The main results which are derived following the abstract formulation in [43] are provided in the next theorem.

**Theorem 3.2.8.** *Let  $\Sigma \triangleq \mathcal{X} \times \mathcal{Y}$ . For every  $\frac{\ell}{s} : \Sigma \rightarrow \mathfrak{R}$  measurable function, bounded below and for  $s > 0$*

$$\begin{aligned} L_4(Q_{X,Y}^*, \lambda^*, s) &= \sup_{Q_{X,Y} \in \mathcal{C}(P_{X,Y})} \left\{ \int_{\Sigma} \ell(x, \Phi(y)) dQ_{X,Y}(x, y) - sH(Q_{X,Y}|P_{X,Y}) \right\} + sR \\ &= s \log \left( \int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y) \right) + sR. \end{aligned} \quad (3.42)$$

Moreover, if  $\ell(x, \Phi(y)) e^{\frac{\ell(x, \Phi(y))}{s}} \in L_1(P_{X,Y})$ , then the supremum in (3.42) is attained by the tilted probability measure  $Q_{X,Y}^*$  given by

$$dQ_{X,Y}^*(x, y) = \frac{e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y)}{\int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y)} \quad (3.43)$$

*Proof.* The derivation is similar to the one presented under Theorem 3.2.5.  $\square$

Next, several properties associated with the maximization over the set  $\mathcal{C}(P_{X,Y})$  are presented.

**Lemma 3.2.9.** *Suppose for a given  $\Phi \in \mathcal{X}_{ad}$ ,*

$$\sup_{Q_{X,Y} \in \mathcal{C}(P_{X,Y})} J_4(\Phi, Q_{X,Y}) < \infty,$$

$s \in \mathfrak{R}$ ,  $\frac{\ell}{s}$  a measurable function bounded from below, the following statements hold.

1) *The dual functional  $L_4(Q_{X,Y}^*, \lambda^*, s)$  is related to the cumulant generating function of  $\ell(x, \Phi(y))$  with respect to  $P_{X,Y} \in \mathcal{M}_1(\mathcal{X} \times \mathcal{Y})$  via*

$$\begin{aligned} L_4(Q_{X,Y}^*, \lambda^*, s) &= s \sup_{\{Q_{X,Y} \in \mathcal{M}_1(\Sigma); H(Q_{X,Y}|P_{X,Y}) < \infty\}} \\ &\quad \left\{ \frac{1}{s} \int_{\Sigma} \ell(x, \Phi(y)) dQ_{X,Y}(x, y) - H(Q_{X,Y}|P_{X,Y}) \right\} + sR \\ &= s \log \int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y) + sR \\ &= s\Psi_{P_{X,Y}}\left(\frac{1}{s}\right) + sR. \end{aligned} \quad (3.44)$$

Moreover, if  $\ell e^{\frac{\ell}{s}} \in L_1(P_{X,Y})$  the supremum in (3.44) is attained at  $Q_{X,Y}^* \in \mathcal{M}(\mathcal{X} \times \mathcal{Y})$  and it is given by

$$dQ_{X,Y}^*(x, y) = \frac{e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y)}{\int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y)}. \quad (3.45)$$

In addition, “The average energy of the system” = “The Helmholtz Free Energy” +  $s \times$  “The Relative Entropy of the system”, that is,

$$\begin{aligned} & \int_{\Sigma} \ell(x, \Phi(y)) dQ_{X,Y}^*(x, y) \\ &= s \log \int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y) + sH(Q_{X,Y}^* | P_{X,Y}), \quad s \in (0, \infty) \end{aligned} \quad (3.46)$$

2) The dual functional  $L_4(Q_{X,Y}^*, \lambda^*, s)$  is convex in  $s > 0$ .

3) The function  $\Gamma_{P_{X,Y}}(s) \triangleq s\Psi_{P_{X,Y}}(\frac{1}{s})$  is a non-increasing function of  $s \in (0, \infty)$ , that is,

$$\Gamma_{P_{X,Y}}(s) = s_1 \log E_{P_{X,Y}} \left\{ e^{\frac{\ell(x, \Phi(y))}{s_1}} \right\} \leq s_2 \log E_{P_{X,Y}} \left\{ e^{\frac{\ell(x, \Phi(y))}{s_2}} \right\} = \Gamma_{P_{X,Y}}(s_2), \quad 0 < s_2 \leq s_1 \quad (3.47)$$

4) The infimum of the dual functional  $L_4(Q_{X,Y}^*, \lambda^*, s)$  over  $s > 0$  defined by

$$L_4(Q_{X,Y}^*, \lambda^*, s^*) = \inf_{s>0} \left\{ s\Psi_{P_{X,Y}}\left(\frac{1}{s}\right) + sR \right\} \quad (3.48)$$

is a concave functional of  $R \geq 0$ .

5) Assume  $\exists \eta > 0$  such that  $\ell e^{\eta \ell} \in L_1(P_{X,Y})$ . Then  $L_4(Q_{X,Y}^*, \lambda^*, s^*)$  evaluated at  $R = 0$  is given by

$$L_4(Q_{X,Y}^*, \lambda^*, s^*) \Big|_{R=0} = \lim_{s \rightarrow \infty} s \log \int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y}(x, y) = E_{P_{X,Y}} \left\{ \ell(x, \Phi(y)) \right\}. \quad (3.49)$$

6) Under the assumptions of 5), the supremum of the dual functional  $L_4(Q_{X,Y}^*, \lambda^*, s)$  over  $s > 0$  is bounded above and from below as follows.

$$E_{P_{X,Y}} \left\{ \ell(x, \Phi(y)) \right\} \leq L_4(Q_{X,Y}^*, \lambda^*, s^*) \leq R + \log E_{P_{X,Y}} \left\{ e^{\ell(x, \Phi(y))} \right\}. \quad (3.50)$$

Moreover if  $\ell(x, \Phi(y))$  is  $Q_{X,Y}$ -essentially bounded for all  $Q_{X,Y} \in \mathcal{C}$ , then the above bounds become

$$\begin{aligned} E_{P_{X,Y}} \left\{ \ell(x, \Phi(y)) \right\} & \leq L_4(Q_{X,Y}^*, \lambda^*, s^*) \\ & \leq \min \left\{ R + \log E_{P_{X,Y}} \left\{ \ell(x, \Phi(y)) \right\}, \|\ell(x, \Phi(y))\|_{\infty} \right\}. \end{aligned}$$

7) If for any  $\eta > 0$ ,  $\ell e^{\eta \ell} \in L_1(P_{X,Y})$  and  $(\ell)^2 e^{\eta \ell} \in L_1(P_{X,Y})$  then the infimum of the functional  $L_4(Q_{X,Y}^*, \lambda^*, s)$  over  $s > 0$  is uniquely attained at

$$H(Q_{X,Y}^* | P_{X,Y})|_{s=s^*} = R \quad (3.51)$$

where  $Q_{X,Y}^*$  is given by (3.45).

That is, a necessary condition for the infimum of the dual functional  $L_4(Q_{X,Y}^*, \lambda^*, s)$  over  $s > 0$  is that  $s^*$  occurs on the boundary of the relative entropy constraint.

Moreover,

$$\begin{aligned} \frac{d}{ds} s \log \int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y} &= \log \int_{\Sigma} e^{\frac{\ell(x, \Phi(y))}{s}} dP_{X,Y} - \frac{1}{s} E_{Q_{X,Y}^*} \{ \ell(x, \Phi(y)) \} \\ &= -H(Q_{X,Y}^* | P_{X,Y}). \end{aligned} \quad (3.52)$$

8) Under the assumptions of 7), the relative entropy  $H(Q_{X,Y}^* | P_{X,Y})$  is a non-increasing function of  $s > 0$ , that is,

$$0 \leq H(Q_{X,Y}^* | P_{X,Y})|_{s=s_2} \leq H(Q_{X,Y}^* | P_{X,Y})|_{s=s_1} \leq H(Q_{X,Y}^* | P_{X,Y})|_{s=s^*} = R, \quad 0 < s^* \leq s_1 \leq s_2. \quad (3.53)$$

*Proof.* The derivations are similar to those in [43] hence they are omitted.  $\square$

It is pointed out that 4) indicates the concavity of dual function with respect to the uncertainty radius, while 5) shows that as  $s \rightarrow \infty$  the estimation problem for a class of distribution will converge to an estimation problem for a single distribution. 6) provides lower and upper bounds on the optimal pay-off as a function of the nominal measure. 8) illustrates the non-increasing property of relative entropy as a function of  $s$ , hence this together with (3.51) gives an algorithm for finding the Lagrange multiplier  $s^*$ .

### 3.3 Examples from Estimation Theory

Next, some examples from estimation theory are presented to illustrate how the results of this chapter apply to practical problems.

### 3.3.1 Estimation of Random Variables

Suppose  $X$  and  $Y$  are RVs defined on  $(\Omega, \mathcal{F}, P)$ , which are related via the nominal model

$$Y = HX + W. \quad (3.54)$$

Hence  $Y : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  the observed RV,  $X : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^n, \mathcal{B}(\mathbb{R}^n))$  the unobserved RV, and  $W : (\Omega, \mathcal{F}) \rightarrow (\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$  the noise RV. Assume  $W$  and  $X$  are independent Gaussian Random Variables,  $N(0, \Sigma_W)$ ,  $\Sigma_W > 0$ ,  $N(0, \Sigma_X)$ ,  $\Sigma_X > 0$ . Moreover, assume (3.54) denotes the nominal model. Then

$$\mu(x, dy) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}} e^{-(y-Hx)^T \frac{\Sigma_W^{-1}}{2} (y-Hx)} dy, \quad (3.55)$$

$$P_X(dx) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_X|^{\frac{1}{2}}} e^{-x^T \frac{\Sigma_X^{-1}}{2} x} dx. \quad (3.56)$$

Let

$$\ell(x, \Phi(y)) = (x - \Phi(y))^T U (x - \Phi(y)), \quad U = U^T > 0. \quad (3.57)$$

Using the above model four different estimation scenarios are investigated. Complete derivations of the results presented below can be found in Appendix C.

#### Application 1

Suppose the uncertainty is described by  $\mathcal{B}^y(\eta) \triangleq \{\nu \in \mathcal{P} : H(\nu|\eta)(y) \leq R(y)\}$ . Then the worst case measure  $\nu^*(y, dx)$  is given by (3.25). Using (3.55), (3.56) and (3.57) this worst case measure is given as

$$\nu^*(y, dx) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma^y(\tilde{s})|^{\frac{1}{2}}} e^{-(x-m^y(\tilde{s}))^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} (x-m^y(\tilde{s}))} dx \quad (3.58)$$

where  $\Sigma^y(\tilde{s}) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{\tilde{s}(y)} \right)^{-1}$  and  $m^y(\tilde{s}) = \Sigma^y(\tilde{s}) \left( H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi(y) \right)$ .

Moreover, using (3.27) together with Bayes' formula  $\eta(y, dx) = \frac{\mu(x, dy) dP_X(x)}{\int_X \mu(x, dy) dP_X(x)}$  the pay-off function is given as

$$L_2(\nu^*, \lambda^*, \tilde{s}^*) = \inf_{\Phi(\cdot)} \inf_{\tilde{s}(\cdot)} \left\{ -\frac{1}{2} \tilde{s}(y) \log \frac{|\Sigma - \frac{2U}{\tilde{s}(y)}|}{|\Sigma|} + \tilde{s}(y) g(\tilde{s}, \Phi, y) + \tilde{s}(y) R(y) \right\} \quad (3.59)$$

where  $\Sigma = H^T \Sigma_W^{-1} H + \Sigma_X^{-1}$  and

$$g(\tilde{s}, \Phi, y) = \left( H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi(y) \right)^T \frac{(\Sigma - \frac{2U}{\tilde{s}(y)})^{-1}}{2} \left( H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi(y) \right) + \Phi(y)^T \frac{U}{\tilde{s}(y)} \Phi(y) - y^T \Sigma_W^{-1} H \frac{\Sigma^{-1}}{2} H^T \Sigma_W^{-1} y.$$

$L_2(\nu^*, \lambda^*, \tilde{s})$  is a convex function of  $\Phi$  for all  $\tilde{s}(\cdot) > 0$  and a convex function of  $\tilde{s}$  for all  $\Phi(\cdot)$ .

In order to calculate the best estimate  $\Phi^*(y)$  of  $X$ , the above pay-off function has to be differentiated and the derivative set to zero as shown below.

$$\frac{d}{d\Phi} L_2(\nu^*, \lambda^*, \tilde{s})|_{\Phi=\Phi^*} = 0, \quad \forall s(\cdot). \quad (3.60)$$

By executing the above differentiation and after some manipulations the best estimate is given by

$$\Phi^*(y) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^T \Sigma_W^{-1} y = \Sigma_X H^T \left( H \Sigma_X H^T + \Sigma_W \right)^{-1} y. \quad (3.61)$$

If it is assumed that the nominal model is the true model, the Least-Square estimate of  $X$  is given by [3]

$$\begin{aligned} \Phi(y) = E[X|Y = y] &= \int_x x \eta(y, dx) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^T \Sigma_W^{-1} y \\ &= \Sigma_X H^T \left( H \Sigma_X H^T + \Sigma_W \right)^{-1} y. \end{aligned} \quad (3.62)$$

This shows that the best estimate  $\Phi^*(y)$  of  $X$  is equal with the expected value of  $X$  given  $Y$ . The reason for this result is the fact that random variables are being used.

Note that when  $\Sigma_W$  is very large compare to  $H \Sigma_X H^T$ , then  $\Phi^*(y) = E[X] = 0$ .

For this problem the minimax problem is equal to the maximin problem therefore, one can perform the minimization of (3.22) with respect to  $\Phi$  to obtain  $\Phi^*(y) = \int x d\nu(y, dx)$ , while the maximization is given by  $\nu^*(y, dx)$  given by (3.25) in which  $\Phi$  is replaced by  $\Phi^*$ . Performing the calculation  $\Phi^*(y) = \int x d\nu^*(y, dx)$  in which both left and right side terms involve the term  $\Phi^*$  and solving for this estimator yields (3.61).

## Application 2

It is assumed that the uncertainty is described by

$$\mathcal{A}(\mu) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{X}} H(\nu|\mu)(x) dP_X(x) \leq \int_{\mathcal{X}} R(x) dP_X(x) \stackrel{\nabla}{=} r_1 \right\}.$$

Then the worst case measure  $\nu^*(x, dy)$  is given by (3.15), where  $\mu(x, dy)$  is given by (3.55). This case can be applied to estimation problems when the received signal is affected by a class of uncertain channels. Next the connection of relative entropy uncertainty to parametric uncertainty is described. Assume  $W$  and  $X$  are independent Gaussian Random Variables,  $N(0, \Sigma_W)$ ,  $\Sigma_W > 0$ ,  $N(0, \Sigma_X)$ ,  $\Sigma_X > 0$ . Suppose the uncertain measure  $\nu(x, dy)$  is induced by the following system.

$$Y = (H + \Delta H)X + W, \quad W \sim N(0; \Sigma_W + \Delta \Sigma_W) \quad (3.63)$$

where  $\Delta H$  denotes the uncertainty matrix to be defined shortly, and  $\Delta \Sigma_W \geq 0$  denotes the noise uncertainty. Suppose  $\mu(x, dy)$  corresponds to the nominal system (3.54) ( $\Delta H = 0$ ,  $\Delta \Sigma_W = 0$ ) and  $\nu(x, dy)$  corresponds to the uncertain system ( $\Delta H \neq 0$ ,  $\Delta \Sigma_W \neq 0$ ). Then the relative entropy between  $\mu(x, dy)$  and  $\nu(x, dy)$  is given by the following expression:

$$\begin{aligned} H(\nu|\mu)(x) &= \int \log \left( \frac{\nu(x, dy)}{\mu(x, dy)} \right) \nu(x, dy) \\ &= \frac{1}{2} \left\{ \log \frac{|\Sigma_W|}{|\Sigma_W + \Delta \Sigma_W|} + \text{tr}((\Sigma_W + \Delta \Sigma_W)(\Sigma_W^{-1} - (\Sigma_W + \Delta \Sigma_W)^{-1})) \right. \\ &\quad \left. + x^T (\Delta H)^T \Sigma_W^{-1} (\Delta H) x \right\}. \end{aligned} \quad (3.64)$$

Taking the average with respect to  $x$ ,

$$\begin{aligned} \int_{\mathcal{X}} H(\nu|\mu)(x) dP_X(x) &= \frac{1}{2} \left\{ \log \frac{|\Sigma_W|}{|\Sigma_W + \Delta \Sigma_W|} \right. \\ &\quad \left. + \text{tr}((\Sigma_W + \Delta \Sigma_W)(\Sigma_W^{-1} - (\Sigma_W + \Delta \Sigma_W)^{-1})) \right. \\ &\quad \left. + \text{tr}(\Sigma_X (\Delta H)^T \Sigma_W^{-1} \Delta H) \right\}. \end{aligned} \quad (3.65)$$

Now, in the special case when  $\Delta \Sigma_W = 0$  the previous expression is equivalent to the weighted norm

$$\|\Delta H\|_N^2 = \frac{1}{2} \text{tr} \left( \Sigma_X^{\frac{1}{2}} (\Delta H)^T \Sigma_W^{-1} \Delta H \Sigma_X^{\frac{1}{2}} \right). \quad (3.66)$$



The conclusion by using the above expression is that the norm uncertainty defined by (3.66) is a special case of the relative entropy uncertainty. This can be seen if the following uncertainty set of channels is defined:

$$\mathcal{N} = \left\{ Z \text{ is a channel characterized by } H + \Delta H, W \sim N(0; \Sigma_W) : \right. \\ \left. \sqrt{2} \|\Delta H\|_N \leq \sqrt{r_1} \right\}.$$

Then the output distribution of every channel in  $\mathcal{N}$  belongs to  $\mathcal{A}(\mu)$ , and in order to solve the estimation problem over the whole class of channels  $\mathcal{N}$  (given the observations  $Y$ ), one can solve the original problem as posed in Section 3.2.2. Notice that  $\mathcal{N} \subseteq \mathcal{A}(\mu)$  because relative entropy uncertainty allows more general models than (3.63) with  $\Delta \Sigma_W = 0$ . For example, the uncertain measure noise  $W$  can be different from the noise of the nominal channel in (3.54).

#### Application 3

Suppose the uncertainty is described by

$$\mathcal{C}(P_{X,Y}) \triangleq \{Q_{X,Y} \in \mathcal{M}(\mathcal{X} \times \mathcal{Y}) : H(Q_{X,Y}|P_{X,Y}) \leq R\}.$$

The worst case measure  $dQ_{X,Y}^*(x, y)$  is given by (3.43). Then using (3.55), (3.56) and (3.57) this worst case measure is given by

$$dQ_{X,Y}^*(x, y) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma(s)|^{\frac{1}{2}}} e^{-(x-\kappa^y(s))^T \frac{\Sigma(s)^{-1}}{2} (x-\kappa^y(s))} \frac{e^{\theta(s, \Phi, y)}}{\int_{\mathcal{Y}} e^{\theta(s, \Phi, y)} dy} dy dx \quad (3.67)$$

where

$$\Sigma(s) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - 2 \frac{U}{s} \right)^{-1} \\ \kappa^y(s) = \Sigma(s) \left( H^T \Sigma_W^{-1} y - \frac{2U}{s} \Phi(y) \right)$$

and

$$\theta(s, \Phi, y) = \kappa^y(s)^T \frac{\Sigma(s)^{-1}}{2} \kappa^y(s) + \Phi(y)^T \frac{U}{s} \Phi(y) - y^T \frac{\Sigma_W^{-1}}{2} y.$$

The average pay-off is given by (3.42). Notice that  $L_4(Q^*, \lambda^*, s)$  is a convex function of  $\Phi$  for all  $s > 0$  and a convex function of  $s$  for all  $\Phi(\cdot)$ . So in order to calculate the best estimate  $\Phi^*(Y)$  of  $X$  the average pay-off (3.42) has to be differentiated and the derivative set to zero as follows:

$$\left. \frac{d}{d\Phi} L_4(Q^*, \lambda^*, s) \right|_{\Phi=\Phi^*} = 0, \quad \forall s. \quad (3.68)$$

Performing the above differentiation, and using Bayes' formula

$dP_{X,Y}(x, y) = \eta(y, dx)dP_Y(y) = \mu(x, dy)dP_X(x)$  the best estimate is given by

$$\Phi^*(y) = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)}. \quad (3.69)$$

After some manipulations the best estimate  $\Phi^*(Y)$  of  $X$  is given by

$$\Phi^*(y) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^T \Sigma_W^{-1} y = \Sigma_X H^T \left( H \Sigma_X H^T + \Sigma_W \right)^{-1} y. \quad (3.70)$$

Notice that the  $\Phi^*(Y)$  is the same as (3.61) which is the expected value of  $X$  given  $Y$ .

#### Application 4

Suppose the uncertainty is described by

$$\mathcal{B}^R(\eta) \triangleq \left\{ \nu \in \mathcal{P} : \int_y H(\nu|\eta)(y) dP_Y(y) \leq \int_{\mathcal{X} \times \mathcal{Y}} R(x, y) \nu(y, dx) dP_Y(y) + \bar{R} \triangleq r_3 \right\}$$

and assume that

$$R(x, y) = (y - \bar{H}x)^T \tilde{U} (y - \bar{H}x), \quad \tilde{U} = \tilde{U}^T > 0. \quad (3.71)$$

Then the worst case measure  $\nu^*(y, dx)$  is given by (3.37). Using (3.55), (3.56), (3.57) and (3.71) this worst case measure is given by

$$\nu^*(y, dx) = \frac{1}{(2\pi)^{\frac{n}{2}} |\tilde{\Sigma}(s)|^{\frac{1}{2}}} e^{-(x - \tilde{m}^y(s))^T \frac{\tilde{\Sigma}(s)^{-1}}{2} (x - \tilde{m}^y(s))} dx \quad (3.72)$$

where

$$\begin{aligned} \tilde{\Sigma}(s) &= \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\bar{H}^T \tilde{U} \bar{H} - \frac{2U}{s} \right)^{-1}, \\ \tilde{m}^y(s) &= \tilde{\Sigma}(s) \left( H^T \Sigma_W^{-1} y + 2\bar{H}^T \tilde{U} y - \frac{2U}{s} \Phi(y) \right). \end{aligned}$$

The average pay-off function is given by (3.39). Like the previous application this pay-off function,  $L_3^R(\nu^*, \lambda^*, s)$ , is a convex function of  $\Phi$  for all  $s > 0$  and a convex function of  $s$  for all  $\Phi(\cdot)$ .

Next, by differentiating the average pay-off (3.39) and setting the derivative to zero the best estimate  $\Phi^*(Y)$  of  $X$  is given by

$$\frac{d}{d\Phi} L_3^R(\nu^*, \lambda^*, s) \Big|_{\Phi=\Phi^*} = 0 \Rightarrow \Phi^*(y) = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx)}. \quad (3.73)$$

Finally, after some manipulations the best estimate  $\Phi^*(Y)$  of  $X$  is given by

$$\Phi^*(y) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\bar{H}^T \tilde{U} \bar{H} \right)^{-1} \left( H^T \Sigma_W^{-1} + 2\bar{H}^T \tilde{U} \right) y. \quad (3.74)$$

**Remark 3.3.1.** *The best estimate  $\Phi^*(Y)$  of  $X$  is given by (3.74). Then,*

i. *if it is assumed that  $\bar{H} = H$ ,*

$$\begin{aligned} \Phi^*(y) &= \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2H^T \tilde{U} H \right)^{-1} \left( H^T \Sigma_W^{-1} + 2H^T \tilde{U} \right) y \\ &= \left( H^T (\Sigma_W^{-1} + 2\tilde{U}) H + \Sigma_X^{-1} \right)^{-1} H^T (\Sigma_W^{-1} + 2\tilde{U}) y; \end{aligned} \quad (3.75)$$

ii. *if it is assumed that  $\bar{H} = I$ ,*

$$\Phi^*(y) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\tilde{U} \right)^{-1} \left( \Sigma_W^{-1} + 2\tilde{U} \right) y. \quad (3.76)$$

### 3.3.2 Estimation of a Sequence of Random Variables

Suppose  $X$  is a RV and  $Y^m = \{Y_0, Y_1, \dots, Y_m\}$  is a sequence of RVs defined on  $(\Omega, \mathcal{F}, P)$ , which are related via the nominal model

$$Y_i = HX + W_i, \quad i = 0, \dots, m \quad (3.77)$$

Hence  $Y^m : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}^{(m+1)d}, \mathcal{B}(\mathfrak{R}^{(m+1)d}))$  the observed sequence of RVs,  $X : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}^n, \mathcal{B}(\mathfrak{R}^n))$  the unobserved RV, and  $W^m : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}^{(m+1)d}, \mathcal{B}(\mathfrak{R}^{(m+1)d}))$  the noise sequence of RVs. Assume  $W_i$  and  $X$  are independent Gaussian Random Variables,  $N(0, \Sigma_W)$ ,  $\Sigma_W > 0$ ,  $N(0, \Sigma_X)$ ,  $\Sigma_X > 0$ . Moreover, assume (3.77) denotes the nominal model. Then

$$\mu(x, dy^m) = \prod_{i=0}^m \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}} e^{-(y_i - Hx)^T \frac{\Sigma_W^{-1}}{2} (y_i - Hx)} dy_i, \quad (3.78)$$

$$P_X(dx) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_X|^{\frac{1}{2}}} e^{-x^T \frac{\Sigma_X^{-1}}{2} x} dx. \quad (3.79)$$

Let

$$\begin{aligned} \ell(x, \Phi(y^m)) &= \sum_{i=0}^{m-1} (x - \Phi_i^*(y^i))^T U (x - \Phi_i^*(y^i)) + (x - \Phi_m(y^m))^T U (x - \Phi_m(y^m)) \\ U &= U^T > 0. \end{aligned} \quad (3.80)$$

Suppose the uncertainty is described by  $\mathcal{B}^y(\eta) \triangleq \{\nu \in \mathcal{P} : H(\nu|\eta)(y) \leq R(y)\}$ . Then the worst case measure  $\nu^*(y^m, dx)$  is given by (3.25). Using (3.78), (3.79) and (3.80) this worst case measure can be expressed as

$$\nu^*(y^m, dx) = \frac{1}{(2\pi)^{\frac{n}{2}} |\widehat{\Sigma}^y(\tilde{s})|^{\frac{1}{2}}} e^{-(x-\widehat{m}^y(\tilde{s}))^T \frac{\widehat{\Sigma}^y(\tilde{s})^{-1}}{2} (x-m^y(\tilde{s}))} dx \quad (3.81)$$

where

$$\begin{aligned} \widehat{\Sigma}^y(\tilde{s}) &= \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2(m+1)U}{\tilde{s}(y)} \right)^{-1}, \\ \widehat{m}^y(\tilde{s}) &= \widehat{\Sigma}^y(\tilde{s}) \left( H^T \Sigma_W^{-1} \sum_{i=0}^m y_i - \frac{2U}{\tilde{s}(y)} \sum_{i=0}^{m-1} \Phi_i^*(y^i) - \frac{2U\Phi_m(y^m)}{\tilde{s}(y)} \right). \end{aligned}$$

Note that  $\widehat{\Sigma}^y(\tilde{s})$  is a valid covariance provided it is positive semidefinite. Hence,

$$\begin{aligned} \widehat{\Sigma}^y(\tilde{s}) &= \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2(m+1)U}{\tilde{s}(y)} \right)^{-1} > 0 \\ &\Rightarrow \left( 2(m+1)U \left( \frac{U^{-1}}{2} H^T \Sigma_W^{-1} H + \frac{U^{-1}}{2(m+1)} \Sigma_X^{-1} - \frac{I}{\tilde{s}(y)} \right) \right)^{-1} > 0 \\ &\Rightarrow \left( 2(m+1)U \left( A - \frac{I}{\tilde{s}(y)} \right) \right)^{-1} > 0, \quad A \triangleq \frac{U^{-1}}{2} H^T \Sigma_W^{-1} H + \frac{U^{-1}}{2(m+1)} \Sigma_X^{-1} \\ &\Rightarrow A - \frac{I}{\tilde{s}(y)} > 0. \end{aligned} \quad (3.82)$$

Given that  $A$  is a square matrix (so it can be written in the form  $A = V\Lambda V^{-1}$ ), then (3.82) can be formulated as

$$V\Lambda V^{-1} - \frac{1}{s} VV^{-1} > 0 \Rightarrow V\left(\Lambda - \frac{1}{s}\right)V^{-1} > 0.$$

Therefore,  $s > \frac{1}{\lambda_{\min}}$ , where  $\lambda_{\min}$  is the smallest eigenvalue of  $A$ .

The average pay-off function  $L_2(\nu^*, \lambda^*, s)$  given by (3.27) is a convex function of  $\Phi_m(y^m)$  for all  $\tilde{s}(\cdot) > 0$  and a convex function of  $\tilde{s}$  for all  $\Phi(\cdot)$ . Similar with previous applications presented, the best estimate  $\Phi_m^*(y^m)$  of  $X$  is calculated by differentiating the average pay-off (3.27) and setting the derivative to zero

$$\left. \frac{d}{d\Phi_m(y^m)} L_2(\nu^*, \lambda^*, \tilde{s}) \right|_{\Phi_m(y^m) = \Phi_m^*(y^m)} = 0 \Rightarrow \Phi_m^*(y^m) = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi_m^*(y^m))}{\tilde{s}}} \eta(y^m, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi_m^*(y^m))}{\tilde{s}}} \eta(y^m, dx)}. \quad (3.83)$$

By further manipulating the above expression, the best estimate is given by

$$\Phi_m^*(y^m) = \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2mU}{\tilde{s}(y)} \right)^{-1} \left( H^T \Sigma_W^{-1} \sum_{i=0}^m y_i - \frac{2U}{\tilde{s}(y)} \sum_{i=0}^{m-1} \Phi_i^*(y^i) \right) \quad (3.84)$$

Clearly, (3.84) is a nonlinear estimator because  $\left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2mU}{\tilde{s}(y)} \right)$  is a function of the observed sequence of RVs  $Y^m$ . When  $\tilde{s}(y) \rightarrow \infty$ , then

$$\Phi_m^*(y^m) = \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^T \Sigma_W^{-1} \sum_{i=0}^m y_i = E[X|Y^m = y^m]. \quad (3.85)$$

Note that, complete derivations of the above results can be found in Appendix C.

## 3.4 Examples from MIMO Communication Systems

In this section the theory developed is applied to Multiple-Input Multiple-Output communication systems. A short overview of the MIMO technology is introduced, and then examples are presented in order to illustrate the applicability of the results to MIMO communication systems.

### 3.4.1 Overview of MIMO Communication Systems

In radio communication, Multiple-Input and Multiple-Output (MIMO) systems employ multiple antennas at both the transmitter and receiver to improve communication performance. It is one of several forms of smart antenna technology. MIMO technology has attracted attention in wireless communications, since it offers significant increases in data throughput and link range without additional bandwidth or transmit power. They achieve higher spectral efficiency and link reliability or diversity.

Wireless MIMO systems are capable of delivering large increases in capacity through utilization of parallel communication channels [48], [49], [42]. Appearing first in a series of information theory articles published by members of Bell Labs, MIMO

systems now constitute a major research area in telecommunications. It is also considered to be one of the technologies that have a chance to resolve the bottlenecks of traffic capacity in the forthcoming broadband wireless Internet access networks. Multiple antennas, both at the transmitter and the receiver, create a matrix channel. The key advantage is the possibility of transmitting over several spatial modes of the channel matrix within the same time-frequency slot at no additional power expenditure. In addition, if the channel matrix is known both at the transmitter (TX) and the receiver (RX), certain spatial modes (singular modes) of the matrix channel can be used to maximize the SNR for every realization of the channel. The singular modes can be used to transport independent data streams (to increase data rate), or one may choose to exploit the top mode (associated with the largest singular value) in order to maximize the spatial diversity advantage.

As mentioned above, Multiple-Input Multiple-Output technology has emerged recently as one of the most significant technologies in modern communication. By using MIMO technology an increase in the system capacity and/or an improvement in the quality of service can be achieved. The key to fully utilize the MIMO capacity relies heavily on the requirement of accurate channel estimation. MIMO channel estimation methods can be classified into three categories: training-based methods, blind methods and semi-blind methods. For pure training-based schemes, a long training is necessary in order to obtain a reliable MIMO channel estimate which reduces the system bandwidth efficiency considerably. Blind methods which do not require any training symbols achieve high system throughput at the expense of high computational complexity. Semi-blind schemes on the other hand require less computational complexity than blind methods and fewer training symbols than training-based methods, making them attractive for practical implementation.

In this section the theory developed in previous sections of this chapter is implemented in order to solve some simple examples for MIMO communication systems.

### 3.4.2 Estimation from MIMO communication Systems

A complex representation of a MIMO communication channel with  $n$  transmitting and  $d$  receiving antennas, is considered

$$Y = HX + W \quad (3.86)$$

where  $Y \in \mathbb{C}^d$  is the received signal,  $X \in \mathbb{C}^n$  is the transmitted signal,  $H \in \mathbb{C}^{d \times n}$  is the channel matrix, and  $W$  is the zero-mean circularly-symmetric Gaussian noise independent of the transmitted signal  $X$ . The covariance matrices of  $X$  and  $W$  are denoted by  $\Sigma_X$  and  $\Sigma_W$ , respectively.

Furthermore, it is assumed that the following are known

$$\mu(x, dy) = \pi^{-d} |\Sigma_W|^{-1} e^{-(y-Hx)^\dagger \Sigma_W^{-1} (y-Hx)} dy, \quad (3.87)$$

$$P_X(dx) = \pi^{-n} |\Sigma_X|^{-1} e^{-x^\dagger \Sigma_X^{-1} x} dx. \quad (3.88)$$

Let

$$\ell(x, \Phi(y)) = (x - \Phi(y))^\dagger U (x - \Phi(y)), \quad U = U^\dagger > 0. \quad (3.89)$$

$\dagger$  denotes complex conjugate transpose,  $\Sigma_W \in \mathbb{C}^{d \times d}$ ,  $\Sigma_W = \Sigma_W^\dagger \geq 0$ ,  $\Sigma_X \in \mathbb{C}^{n \times n}$ ,  $\Sigma_X = \Sigma_X^\dagger \geq 0$ .

Using the above model two different estimation scenarios are being investigated. Complete derivations of the results presented below can be found in Appendix C.

#### Application 1

Suppose the uncertainty is described by

$$\mathcal{B}^R(\eta) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{Y}} H(\nu|\eta)(y) dP_Y(y) \leq \int_{\mathcal{X} \times \mathcal{Y}} R(x, y) \nu(y, dx) dP_Y(y) + \bar{R} \stackrel{\nabla}{=} r_3 \right\}$$

and assume that

$$R(x, y) = (y - \bar{H}x)^\dagger \tilde{U} (y - \bar{H}x), \quad \tilde{U} = \tilde{U}^\dagger > 0. \quad (3.90)$$

Then the worst case measure  $\nu^*(y, dx)$  is given by (3.37). Using (3.87), (3.88), (3.89) and (3.90) this worst case measure can be expressed as

$$\nu^*(y, dx) = \pi^{-n} |\tilde{\Sigma}_M(s)|^{-1} e^{-(x - \tilde{m}_M(s))^\dagger \tilde{\Sigma}_M(s)^{-1} (x - \tilde{m}_M(s))} dx \quad (3.91)$$

where

$$\begin{aligned} \tilde{\Sigma}_M(s) &= \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \bar{H}^\dagger \tilde{U} \bar{H} - \frac{U}{s} \right)^{-1} \\ \tilde{m}_M(s) &= \tilde{\Sigma}(s) \left( H^\dagger \Sigma_W^{-1} y + \bar{H}^\dagger \tilde{U} y - \frac{U}{s} \Phi(y) \right). \end{aligned}$$

Notice that the pay-off function  $L_3^R(\nu^*, \lambda^*, s)$ , given by (3.39), is a convex function of  $\Phi$  for all  $s > 0$  and a convex function of  $s$  for all  $\Phi(\cdot)$ . Therefore, by differentiating the average pay-off (3.39) and setting the derivative to zero the best estimate  $\Phi^*(Y)$  of  $X$  is given by

$$\frac{d}{d\Phi} L_3^R(\nu^*, \lambda^*, s) \Big|_{\Phi=\Phi^*} = 0 \Rightarrow \Phi^*(y) = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y)) - R(x, y)}{s}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y)) - R(x, y)}{s}} \eta(y, dx)}. \quad (3.92)$$

Finally, after some manipulations of the above expression, the best estimate is given by

$$\Phi^*(y) = \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \bar{H}^\dagger \tilde{U} \bar{H} \right)^{-1} \left( H^\dagger \Sigma_W^{-1} + \bar{H}^\dagger \tilde{U} \right) y. \quad (3.93)$$

**Remark 3.4.1.** The best estimate  $\Phi^*(Y)$  of  $X$  is given by (3.93). Then,

i. if it is assumed that  $\bar{H} = H$ ,

$$\begin{aligned} \Phi^*(y) &= \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + H^\dagger \tilde{U} H \right)^{-1} \left( H^\dagger \Sigma_W^{-1} + H^\dagger \tilde{U} \right) y \\ &= \left( H^\dagger (\Sigma_W^{-1} + \tilde{U}) H + \Sigma_X^{-1} \right)^{-1} H^\dagger (\Sigma_W^{-1} + \tilde{U}) y; \end{aligned} \quad (3.94)$$

ii. if it is assumed that  $\bar{H} = I$ ,

$$\Phi^*(y) = \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \tilde{U} \right)^{-1} \left( H^\dagger \Sigma_W^{-1} + \tilde{U} \right) y. \quad (3.95)$$

### Application 2

Suppose the uncertainty is described by

$$\mathcal{C}(P_{X,Y}) \triangleq \{Q_{X,Y} \in \mathcal{M}(\mathcal{X} \times \mathcal{Y}) : H(Q_{X,Y} | P_{X,Y}) \leq R\}.$$

The worst case measure  $Q_{X,Y}^*(x, y)$  is given by (3.43). Then using (3.87), (3.88) and (3.89) this worst case measure is given by

$$dQ_{X,Y}^*(x, y) = (2\pi)^{-n} |\Sigma_M(s)|^{-1} e^{-(x - \kappa_M^y(s))^\dagger \Sigma_M(s)^{-1} (x - \kappa_M^y(s))} \frac{e^{\tilde{\theta}(s, \Phi, y)}}{\int_{\mathcal{Y}} e^{\tilde{\theta}(s, \Phi, y)} dy} dy dx \quad (3.96)$$

where

$$\begin{aligned} \Sigma_M(s) &= \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{U}{s} \right)^{-1} \\ \kappa_M^y(s) &= \Sigma_M(s) \left( H^\dagger \Sigma_W^{-1} y - \frac{U}{s} \Phi(y) \right) \end{aligned}$$



and

$$\tilde{\theta}(s, \Phi, y) = \kappa_M^y(s)^\dagger \Sigma_M(s)^{-1} \kappa_M^y(s) + \Phi(y)^\dagger \frac{U}{s} \Phi(y) - y^\dagger \Sigma_W^{-1} y.$$

Given that the average pay-off  $L_4(\nu^*, \lambda^*, s)$ , given by (3.42) is a convex function of  $\Phi$  for all  $s > 0$  and a convex function of  $s$  for all  $\Phi(\cdot)$ , in order to calculate the best estimate  $\Phi^*(Y)$  of  $X$ , one just has to differentiate the average pay-off (3.42) and set the derivative to zero as follows

$$\left. \frac{d}{d\Phi} L_4(\nu^*, \lambda^*, s) \right|_{\Phi=\Phi^*} = 0 \Rightarrow \Phi^*(y) = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)}. \quad (3.97)$$

Then, after some calculations and manipulations the best estimate  $\Phi^*(Y)$  of  $X$  is given by

$$\Phi^*(y) = \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^\dagger \Sigma_W^{-1} y = \Sigma_X H^\dagger \left( H \Sigma_X H^\dagger + \Sigma_W \right)^{-1} y. \quad (3.98)$$

### 3.5 Summary

This chapter considers the problem of least-square estimation for a class of systems which are subject to uncertainty, and employs the KL distance to describe the uncertainty classes. Stochastic kernels, and joint distributions are used to describe the uncertainty models and a minimax approach is implemented. Three problems are formulated and their solutions are sought, highlighting some properties associated with the estimate of the true distribution. Classical examples are chosen to illustrate the applicability of the results. The theory developed is also applied to simple examples for MIMO communication systems.



# CHAPTER 4

## APPLICATIONS OF ROBUST ESTIMATION

In this chapter, the theory developed in Chapter 3 is applied to a nonlinear recursive model in order to derive robust estimators. The methodology presented invokes a change of probability measure technique to derive recursive equations for the conditional density of nonlinear filtering problems. In section 4.2, the mathematical theory is presented by extending the theory of Chapter 3. In Section 4.3, the theory developed in Section 4.2 is applied to a linear Gaussian example. While in Section 4.4, the theory is applied to an attenuated sinusoid in a multipath environment, where various estimators are derived by solving the recursive equation satisfied by the unnormalized version of the a posteriori density.

### 4.1 Introduction

Chapter 4 deals with estimation techniques, for estimating signals which are generated by finite-dimensional autoregressive channel models found in [34], [35], subject to uncertainty. Autoregressive channel models have been used with success to predict fading channel dynamics for the purposes of Kalman filter based channel estimation and for long-range channel forecasting. They have also been used by several authors to simulate correlated Rayleigh fading. The uncertainty description of these systems is characterized by the class of uncertain measures which satisfy a Kullback-Leibler (KL) distance constraint with respect to a nominal measure. The uncertainty description of the system, and the nominal description of the system are modeled by joint probability distributions.

A minimax approach is implemented here in order to address the above problem. The difficulty with this approach lies in the minimization. Because this is a nonlinear problem it is not easy to obtain the solution of the  $\hat{a}$  posteriori distribution recursion in closed form. Very often, when dealing with nonlinear estimation problems, approximations have to be made to find sub-optimal nonlinear estimators. A common method is to use the so-called extended Kalman filter [5, 6, 7]. Other more sophisticated sub-optimal estimation techniques are available, e.g., reiteration, higher order filters, and statistical linearization [5]. The chapter deals specifically with the problem of finding a solution to the minimax formulation and does not take into consideration other nonlinear estimation techniques.

In general, minimax estimation techniques lead to strategies, which are robust with respect to variations in the models as long as these belong to the uncertainty class. As was already described in Chapter 1 Minimax Wiener filtering techniques are given in [30], while blind minimax linear regression problems of estimating deterministic parameters are given in [23]. Relations of blind minimax techniques to Stein type estimators [27] and least-square regularization [9] is discussed in [23]. Related work in which uncertainty is described by relative entropy can be found in [16], [29], [50]. The examples provided in [16], [23], [29] and [50] are linear Gaussian.

In this chapter, in the abstract setting, the maximization is addressed using variational methods, while the minimization is addressed using a change of probability measure technique [36]. The change of probability measure techniques introduced is being used in order to derive a recursive equation for the unnormalized  $\hat{a}$  posteriori distribution of nonlinear filtering problems. The theory developed is applied first to a linear Gaussian model and then to an attenuated sinusoid in a multipath environment, which is subject to an additive Gaussian noise. The multipath model employed for non-coherent estimation and detection compliments the work found in the literature [4], [51], in the sense that the classical problem assumes no multipath scenario, while the attenuation of the sinusoidal is assumed to be a known deterministic function. Note that, unlike the classical non-coherent estimation problem [1], [4], which is concerned with a fixed model, this chapter deals with a class of models while the estimation is formulated using minimax techniques. The connection to least-squares estimation found in [1], [4] for single channel, is established by reducing the uncertainty to zero, while generalizing existing results.

## 4.2 The Minimax Filtering

In this section, the nonlinear model is described. The minimax problem is defined and a change of probability measure technique is introduced, which reformulates the problem under a new fictitious probability measure, where the signal to be estimated and the observations are independent. Finally, a linear recursion for an unnormalized conditional density related to the minimax filtering is derived.

### 4.2.1 State and Observation Models

Let  $(\Omega, \mathcal{F}, P)$  be complete probability space on which the nominal state or unobserved process  $\{x_k\}, k \in N_0 \triangleq \{0, 1, 2, 3, \dots\}$  and the observation process  $\{y_k\}, k \in N_0$ , are defined by the following recursions:

$$\begin{aligned} x_{k+1} &= f(k+1, x_k) + B_{k+1}w_{k+1}, \quad x_0 \in \mathfrak{R}^n \\ y_k &= h(k, x_k) + D_kv_k, \quad y_0 \in \mathfrak{R}^d. \end{aligned} \quad (4.1)$$

Here  $x_0 : \Omega \rightarrow \mathfrak{R}^n$  is the initial state and  $w : \Omega \times N_0 \rightarrow \mathfrak{R}^n$ ,  $v : \Omega \times N_0 \rightarrow \mathfrak{R}^d$ , are random noises.

Also  $\{w_k\}, \{v_k\}, k \in N_0$  are independent noise sequences of Random Variables (RVs) with densities  $\Psi_w(w) = \frac{1}{(2\pi)^{\frac{n}{2}}} e^{-\frac{w^T w}{2}}$ ,  $\Xi_v(v) = \frac{1}{(2\pi)^{\frac{d}{2}}} e^{-\frac{v^T v}{2}}$ , respectively,  $x_0$  has density  $\pi_{x_0}(x) = \frac{d\Pi_{x_0}(x)}{dx}$ , which is also independent of  $\{w_k\}, \{v_k\}$ .

It is assumed throughout this chapter that  $f : N_0 \times \mathfrak{R}^n \rightarrow \mathfrak{R}^n$  and  $h : N_0 \times \mathfrak{R}^n \rightarrow \mathfrak{R}^d$ ,  $B_k, D_k$  are Borel measurable functions, and  $(B_k B_k^T)^{-1}, (D_k D_k^T)^{-1}$  exist.

### 4.2.2 Definition of Minimax Problem

Let  $\{\mathcal{G}_m^0\}$  be the  $\sigma$ -field generated by the complete data  $\{x_0, x_1, \dots, x_m, y_0, y_1, \dots, y_m\}$  and let  $\{\mathcal{G}_m\}, m \in N_0$  denote its complete filtration [36]. Let  $\{\mathcal{F}_m^{0, \mathcal{Y}}\}$  define the  $\sigma$ -field generated by the incomplete data  $\{y_0, y_1, \dots, y_m\}$  and let  $\{\mathcal{F}_m^{\mathcal{Y}}\}, m \in N_0$  denote its complete filtration. Let  $y^m$  denote the sequence  $\{y_0, \dots, y_m\}$  and similarly for other sequences.  $\tilde{x}_m$  denotes the estimate of the state  $x_m$  given  $\{\mathcal{F}_m^{\mathcal{Y}}\}, m \in N_0$ ; it is assumed recursive estimates which update  $\tilde{x}_m$  from knowledge of  $\tilde{x}_{m-1}$  and

past and present data  $\{y_0, \dots, y_m\}$ ,  $m \in N_0$ . Throughout the chapter  $E_Q[\cdot]$  denotes expectation with respect to probability distribution  $Q$ .

In the next definition the class of admissible estimators, which are functions of the observation sequences are introduced.

**Definition 4.2.1.** *The set of admissible estimators  $\mathcal{X}_{ad}$  are defined as follows.*

$\mathcal{X}_{ad} \triangleq \left\{ \tilde{x} : \Omega \times N_0 \rightarrow X \subseteq \mathbb{R}^n; \{\tilde{x}_k\}, \text{ is adapted to } \{\mathcal{F}_k^y\}, k \in N_0 \right\}$  (e.g., adapted means that at each  $k \in N_0$ ,  $\tilde{x}_k$  is a causal function of the data  $\{\mathcal{F}_k^y\}$ ).

Let  $P_{x^m, y^m}$  denote the nominal (in the absence of modelling uncertainties) joint distribution of the sequences  $(x^m, y^m)$ , which corresponds to the one induced by model (4.1). Let  $Q_{x^m, y^m}$  denote the true joint distribution of the sequences  $(x^m, y^m)$ , which is unknown. The only available information is that  $Q_{x^m, y^m}$  belongs to a class of possible distributions. This class is modelled by the information theoretic relative entropy set

$$\mathcal{C}(P_{x^m, y^m}) \triangleq \{Q_{x^m, y^m} : H(Q_{x^m, y^m} | P_{x^m, y^m}) \leq R\} \quad (4.2)$$

where  $R > 0$  and  $H(\cdot | \cdot)$  is the KL distance between the two joint distributions defined in Section 3.2.4. Note that this is the same uncertainty set used in Section 3.2.4.

The pay-off is the average of a function of the error over the time horizon  $[0, m]$ ,  $E_{Q_{x^m, y^m}} \left\{ \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right\}$ , in which the average is taken with the unknown distribution  $Q_{x^m, y^m} \in \mathcal{C}(P_{x^m, y^m})$ . Here  $\tilde{\ell}(x_k, \tilde{x}_k)$  is a measure of distance between the state  $x_k$  and its estimate  $\tilde{x}_k$  (e.g.  $\|x_k - \tilde{x}_k\|_{\mathbb{R}^n}^2$ ).

Next, the minimax estimation problem is defined.

**Problem 4.2.1.** *Given the nominal probability distribution  $P_{x^m, y^m}$  induced by system (4.1), find a probability distribution  $Q_{x^m, y^m}^*$  and an estimator  $\tilde{x}_m^* \in \mathcal{X}_{ad}$  which solve*

$$J(\tilde{x}^*, Q_{x^m, y^m}^*) = \inf_{\tilde{x}_m \in \mathcal{X}_{ad}} \sup_{Q_{x^m, y^m} \in \mathcal{C}(P_{x^m, y^m})} E_{Q_{x^m, y^m}} \left\{ \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right\} \quad (4.3)$$

when  $R \in (0, \infty)$ , and  $\tilde{\ell} : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, \infty)$ , is continuous in  $(x, \tilde{x}) \in \mathbb{R}^n \times \mathbb{R}^n$  and bounded from below.

The assumptions on  $\tilde{\ell}$  are sufficient for the existence of estimates of the maximizing distribution  $Q_{x^m, y^m}^*$  [43].

### 4.2.3 Minimax Optimization

Define the sample pay-off  $\ell(x, \tilde{x}) \triangleq \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k)$ . For a given  $\tilde{x}_m \in \mathcal{X}_{ad}$ , and a fixed nominal probability distribution  $P_{x^m, y^m}$  induced by the nominal system (4.1), let  $Q_{x^m, y^m}^*$  denote the distribution which achieves the supremum of the average energy pay-off functional subject to the relative entropy uncertainty as defined by

$$J(\tilde{x}_m, Q_{x^m, y^m}^*) = \sup_{Q_{x^m, y^m} \in \mathcal{C}(P_{x^m, y^m})} \int_{\mathfrak{R}^{(m+1)n} \times \mathfrak{R}^{(m+1)d}} \ell(x^m, \tilde{x}^m) dQ_{x^m, y^m}(x, y). \quad (4.4)$$

The solution to this optimization problem returns the worst case probability distribution among those which satisfy the constraint, as a function of the nominal distribution  $P_{x^m, y^m}$ .

Thus, the estimation problem for the class of models  $\mathcal{C}(P_{x^m, y^m})$ , is to find an estimator  $\tilde{x}_m^* \in \mathcal{X}_{ad}$  which solves

$$J(\tilde{x}_m^*) = \inf_{\tilde{x}_m \in \mathcal{X}_{ad}} J(\tilde{x}_m, Q_{x^m, y^m}^*). \quad (4.5)$$

Since (4.4) is a constraint optimization, to find the supremum over  $Q_{x^m, y^m} \in \mathcal{C}(P_{x^m, y^m})$  the Lagrangian has to be defined [40]

$$L(Q_{x^m, y^m}^*, \tilde{x}_m, \lambda^*, s^*) \triangleq \inf_{s \geq 0} \inf_{\lambda \geq 0} \sup_{Q_{x^m, y^m} \in \mathcal{C}(P_{x^m, y^m})} \left\{ \int_{\mathfrak{R}^{(m+1)n} \times \mathfrak{R}^{(m+1)d}} \ell(x^m, \tilde{x}^m) dQ_{x^m, y^m}(x, y) - s \left( H(Q_{x^m, y^m} | P_{x^m, y^m}) - R \right) - \lambda \left( \int_{\mathfrak{R}^{(m+1)n} \times \mathfrak{R}^{(m+1)d}} dQ_{x^m, y^m}(x, y) - 1 \right) \right\} \quad (4.6)$$

where  $s \geq 0$  is the Lagrange multiplier associated with the constraint  $\mathcal{C}(P_{x^m, y^m})$  and  $\lambda \geq 0$  is the Lagrange multiplier associated with the constraint

$\int_{\mathfrak{R}^{(m+1)n} \times \mathfrak{R}^{(m+1)d}} dQ_{x^m, y^m}(x, y) = 1$  (e.g.,  $Q_{x^m, y^m}$  is a joint distribution and hence it should integrate to 1).

From Section 3.2.4 the worst case measure,  $Q_{x^m, y^m}^*$ , is given by

$$dQ_{x^m, y^m}^*(x, y) = \frac{e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} dP_{x^m, y^m}(x, y)}{\int_{\mathfrak{R}^{(m+1)n} \times \mathfrak{R}^{(m+1)d}} e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} dP_{x^m, y^m}(x, y)}. \quad (4.7)$$

Then (4.7) is substituted into (4.6) to deduce

$$L(Q_{x^m, y^m}^*, \tilde{x}_m, \lambda^*, s) = s \log \int_{\mathfrak{R}^{(m+1)n} \times \mathfrak{R}^{(m+1)d}} e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} dP_{x^m, y^m}(x, y) + sR. \quad (4.8)$$

Hence, Problem 4.2.1 is equivalent to the dual problem of finding the optimal estimator  $\tilde{x}_m^* \in \mathcal{X}_{ad}$  via

$$\inf_{s \geq 0} \inf_{\tilde{x}_m \in \mathcal{X}_{ad}} L(Q_{x^m, y^m}^*, \tilde{x}_m, \lambda^*, s). \quad (4.9)$$

Thus, the rest of this chapter deals with the problem of finding the minimum over  $\tilde{x}_m \in \mathcal{X}_{ad}$  of (4.8).

This is achieved by expressing the minimization over  $\tilde{x}_m \in \mathcal{X}_{ad}$  in terms of the solution of a recursive equation satisfied by a conditional unnormalized distribution. The next theorem presents intermediate steps using the change of probability measure discussed extensively in Section 2.2.

**Theorem 4.2.2.** *Let  $\Phi$  be a bounded continuous function on  $\mathfrak{R}^n$  taking values in  $\mathfrak{R}$ . The likelihood function of the complete data  $\{x_0, \dots, x_m, y_0, \dots, y_m\}$  is defined by*

$$\Lambda_m \triangleq \prod_{k=0}^m \left[ \frac{\Xi_v(D_k^{-1}(y_k - h(k, x_k))) \Psi_w(B_k^{-1}(x_k - f(k, x_{k-1})))}{|D_k| \Xi_v(y_k) |B_k| \Psi_w(x_k)} \right] = \frac{dP(x^m, y^m)}{d\bar{P}(x^m, y^m)} \quad (4.10)$$

where under probability distribution  $\bar{P}$ ,  $\{x_k\}$  is i.i.d.  $N(0, I_n)$  and  $\{y_k\}$  is i.i.d.  $N(0, I_d)$  with density functions  $\Psi_w(x_k) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp(\frac{-x_k^T x_k}{2})$  and  $\Xi_v(y_k) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp(\frac{-y_k^T y_k}{2})$ ,  $k \in N_0$ , respectively. Then the following relations hold.

1) Conditional expectations are related via

$$E \left[ \Phi(x_m) \exp \left( \frac{1}{s} \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right) \middle| \mathcal{F}_m^y \right] = \frac{\bar{E} \left[ \Phi(x_m) \Lambda_m \exp \left( \frac{1}{s} \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right) \middle| \mathcal{F}_m^y \right]}{\bar{E} \left[ \Lambda_m \middle| \mathcal{F}_m^y \right]} \quad (4.11)$$

where  $E, \bar{E}$  denotes expectation under the probability distribution  $P, \bar{P}$ , respectively.

Moreover, the numerator of (4.11) can be written in terms of unnormalized probability distribution  $\alpha_m^s(\cdot)$  via

$$\alpha_m^s(\Phi) \triangleq \bar{E} \left[ \Phi(x_m) \Lambda_m \exp \left( \frac{1}{s} \sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k) \right) \middle| \mathcal{F}_m^y \right] = \int_{\mathfrak{R}^n} \Phi(z) d\alpha_m^s(z) \quad (4.12)$$



where the probability distribution  $\alpha_m^s(\cdot)$  satisfies the following recursion

$$\begin{aligned} \alpha_m^s(\Phi) = & \int_{\mathfrak{R}^n} \Phi(x) \exp\left(\frac{1}{s}\tilde{\ell}(x, \tilde{x}_m)\right) \frac{\Xi_v(D_m^{-1}(y_m - h(m, x)))}{|D_m|\Xi_v(y_m)} \\ & \times \alpha_{m-1}^s\left(\frac{\Psi_w(B_m^{-1}(x - f(m, \cdot)))}{|B_m|}\right) dx \end{aligned} \quad (4.13)$$

with initial condition

$$\alpha_0^s(\Phi) = \int_{\mathfrak{R}^n} \Phi(x_0) \exp\left(\frac{1}{s}\tilde{\ell}(x_0, \tilde{x}_0)\right) \frac{\Xi_v\left(D_0^{-1}(y_0 - h(0, x_0))\right)}{|D_0|\Xi_v(y_0)} d\Pi_{x_0}(x_0). \quad (4.14)$$

2) If  $\alpha_m^s(\Phi)$  has a density, that is  $\frac{d}{dx}\alpha_m^s(x) = \bar{\alpha}_m^s(x)$ , then

$$\alpha_m^s(\Phi) = \int_{\mathfrak{R}^n} \Phi(x) d\alpha_m^s(x) = \int_{\mathfrak{R}^n} \Phi(x) \bar{\alpha}_m^s(x) dx \quad (4.15)$$

and the density  $\{\bar{\alpha}_m^s(x)\}_{m \geq 0}$  satisfies the recursion

$$\begin{aligned} \bar{\alpha}_m^s(x) = & \frac{\Xi_v(D_m^{-1}(y_m - h(m, x)))}{|D_m|\Xi_v(y_m)} \exp\left(\frac{1}{s}\tilde{\ell}(x, \tilde{x}_m)\right) \\ & \times \int_{\mathfrak{R}^n} \frac{\Psi_w(B_m^{-1}(x - f(m, z)))}{|B_m|} \bar{\alpha}_{m-1}^s(z) dz \end{aligned} \quad (4.16)$$

with initial condition

$$\bar{\alpha}_0^s(x) = \frac{\Xi_v(D_0^{-1}(y_0 - h(0, x_0)))}{|D_0|\Xi_v(y_0)} \exp\left(\frac{1}{s}\tilde{\ell}(x_0, \tilde{x}_0)\right) \pi_{x_0}(x_0). \quad (4.17)$$

*Proof.* The derivation is a variant of the one found in Theorem 2.2.19 in Section 2.2.4.

1) This follows from the relation of conditional expectation under different probability distribution [43]. (4.12) follows from Theorem 2.2.19 in Section 2.2.4 by absorbing the exponential term  $\exp\left(\frac{1}{s}\sum_{k=0}^m \tilde{\ell}(x_k, \tilde{x}_k)\right)$  into the likelihood function  $\Lambda_m$ .

2) Follows from Theorem 2.2.19 in Section 2.2.4 and the above discussion.

□

Notice that  $\{\bar{\alpha}_m^s(x); m \in N_0\}$  is a sufficient statistic for the robust (e.g., minimax) estimation problem because the initial optimization Problem 4.2.1 is now expressed

in terms of this quantity. Using Theorem 4.2.2, the dual functional (4.8) is now given by

$$L(Q_{x^m, y^m}^*, \tilde{x}_m, \lambda^*, s) = s \log E[e^{\frac{\ell(x^m, \tilde{x}^m)}{s}}] + sR = s \log \overline{E}[e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} \Lambda_m] + sR \quad (4.18)$$

$$= s \log \overline{E} \left[ \overline{E} \left[ e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} \Lambda_m | \mathcal{F}_m^{\mathcal{Y}} \right] \right] + sR \quad (4.19)$$

$$= s \log \int_{\mathfrak{R}^{(m+1)d}} \overline{E} \left[ e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} \Lambda_m | \mathcal{F}_m^{\mathcal{Y}} \right] dP_{y^m}(y) + sR \quad (4.20)$$

$$= s \log \int_{\mathfrak{R}^{(m+1)d}} \alpha_m^s(1) dP_{y^m}(y) + sR. \quad (4.21)$$

The chain of equalities in (4.18)-(4.21) are obtained as follows. Clearly, (4.18) follows by substituting into  $\int e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} dP_{x^m, y^m}$  the density  $dP = \Lambda_m d\overline{P}$ . Also, (4.19) is a property of conditional expectation. In addition, (4.20) follows from (4.19) since  $\overline{E}[e^{\frac{\ell(x^m, \tilde{x}^m)}{s}} \Lambda_m | \mathcal{F}_m^{\mathcal{Y}}] = f(y^m)$ . Finally, (4.21) follows from (4.12) by letting  $\Phi = 1$ .

Further, to find the optimal estimator  $\tilde{x}^*$  one needs to solve the recursive equation (4.13), then substitute into (4.21) and perform the minimization over  $\tilde{x}_m \in \mathcal{X}_{ad}$ . The point to be made in this section is that the minimax estimation yields the estimation problem (4.8) with an exponential cost, while (4.8) is further expressed in terms of the unnormalized conditional distribution  $\{\alpha_m^s(x) : m \in N_0\}$ .

**Remark 4.2.3.** *In decision applications, there will be one distribution  $\alpha_m^s(x)$  for each hypothesis, e.g., in binary hypothesis the decision rule takes the form*

$$\frac{\inf_{\tilde{x}_m \in \mathcal{X}_{ad}} \int_{\mathfrak{R}^n} d\alpha_m^{s,1}(x)}{\inf_{\tilde{x}_m \in \mathcal{X}_{ad}} \int_{\mathfrak{R}^n} d\alpha_m^{s,2}(x)} \underset{H_2}{\overset{H_1}{\geq}} \gamma$$

where  $\gamma$  is the threshold. Clearly,  $\alpha_m^s(x)$  is a sufficient statistic for estimation and decision problems, since both the least-square estimation and decision rule are constructed from this quantity.

In the next sections the results obtained in this section are applied to different problems where model (4.1) has a specific form.

### 4.3 Minimax Estimation for Linear Gaussian Models

The partially observed Gaussian version of (4.1) given by the following recursions, is being considered in this section.

$$(\Omega, \mathcal{F}, \{\mathcal{G}_k\}, P) : \begin{cases} x_{k+1} = A_{k+1}x_k + B_{k+1}w_{k+1}, & x_k \in \mathfrak{R}^n \\ y_k = C_kx_k + D_kv_k, & y_k \in \mathfrak{R}^d. \end{cases} \quad (4.22)$$

The noises are Gaussian distributed as follows,  $w_k \sim N(0, I_n)$ ,  $v_k \sim N(0, I_d)$ , while sequences  $\{w_k\}$ ,  $\{v_k\}$  are mutually independent and independent of  $x_0 \sim N(\bar{x}_0, V_0)$

According to Theorem 4.2.2 equation (4.16),  $\{\bar{\alpha}_m^s(x)\}$  satisfies the following recursion

$$\begin{aligned} \bar{\alpha}_m^s(x) &= \frac{\Xi_v(D_m^{-1}(y_m - C_mx))}{|D_m|\Xi_v(y_m)} \exp\left(\frac{1}{s}\tilde{\ell}(x, \tilde{x}_m)\right) \\ &\times \int_{\mathfrak{R}^n} \frac{\Psi_w(B_m^{-1}(x - A_mz))}{|B_m|} \bar{\alpha}_{m-1}(z) dz \end{aligned} \quad (4.23)$$

where

$$\Psi_w(x) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left(-\frac{x^T x}{2}\right), \quad (4.24)$$

$$\Xi_v(y) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{y^T y}{2}\right), \quad (4.25)$$

$$\ell(x, \tilde{x}_m) = (x - \tilde{x}_m)^T \frac{W_m}{2} (x - \tilde{x}_m). \quad (4.26)$$

A solution of (4.23) is assumed to have the following form

$$\bar{\alpha}_m^s(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |V_{m|m}^s|^{\frac{1}{2}}} \exp\left(- (x - \hat{x}_{m|m}^s)^T \frac{(V_{m|m}^s)^{-1}}{2} (x - \hat{x}_{m|m}^s) + \beta_{m|m}^s\right) \quad (4.27)$$

where  $\{V_{m|m}^s, \hat{x}_{m|m}^s, \beta_{m|m}^s\}$  will be identified shortly. By substituting (4.27) into the recursion (4.23) and using (4.24), (4.25), (4.26) the following recursive relations for  $\{V_{m|m}^s, \hat{x}_{m|m}^s, \beta_{m|m}^s\}$  are deduced.

$$\hat{x}_{m|m}^s = V_{m|m}^s \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1}^s)^{-1} \hat{x}_{m|m-1}^s - \frac{W_m}{s} \tilde{x}_m \right] \quad (4.28)$$

$$V_{m|m}^s = \left( C_m^T (D_m D_m^T)^{-1} C_m + (V_{m|m-1}^s)^{-1} - \frac{W_m}{s} \right)^{-1} \quad (4.29)$$

$$\beta_{m|m}^s = (\hat{x}_{m|m}^s)^T \frac{(V_{m|m}^s)^{-1}}{2} \hat{x}_{m|m}^s - (\hat{x}_{m|m-1}^s)^T \frac{(V_{m|m-1}^s)^{-1}}{2} \hat{x}_{m|m-1}^s + \tilde{x}_m^T \frac{W_m}{2s} \tilde{x}_m$$

$$\begin{aligned}
 & -\frac{1}{2} \log \left[ |V_{m-1|m-1}^s| |A_m^T (B_m B_m^T)^{-1} A + (V_{m-1|m-1}^s)^{-1}| |(V_{m|m}^s)^{-1}| \right] \\
 & - \log [|D_m B_m|] + \frac{1}{2} y_m^T \left( I_d - (D_m D_m^T)^{-1} \right) y_m + \beta_{m-1|m-1}^s
 \end{aligned} \tag{4.30}$$

where

$$\hat{x}_{m|m-1}^s = A_m \hat{x}_{m-1|m-1}^s \tag{4.31}$$

$$V_{m|m-1}^s = B_m B_m^T + A_m V_{m-1|m-1}^s A_m^T \tag{4.32}$$

with initial conditions

$$\hat{x}_{0|0}^s = V_{0|0}^s \left[ C_0^T (D_0^T D_0)^{-1} y_0 - \frac{W_0}{s} \tilde{x}_0 + V_0^{-1} \bar{x}_0 \right] \tag{4.33}$$

$$V_{0|0}^s = \left( C_0^T (D_0^T D_0)^{-1} C_0 - \frac{W_0}{s} + V_0^{-1} \right)^{-1} \tag{4.34}$$

$$\begin{aligned}
 \beta_{0|0}^s &= (\hat{x}_{0|0}^s)^T \frac{(V_{0|0}^s)^{-1}}{2} \hat{x}_{0|0}^s + \tilde{x}_0^T \frac{W_0}{2s} \tilde{x}_0 - \bar{x}_0 \frac{V_0^{-1}}{2} \bar{x}_0 \\
 &+ \frac{1}{2} \log \left[ \frac{V_{0|0}}{V_0} \right] - \log |D_0| |B_0| + \frac{1}{2} y_0^T \left( I_d - (D_0^T D_0)^{-1} \right) y_0
 \end{aligned} \tag{4.35}$$

thus it is concluded that (4.27) is indeed a solution of the recursion (4.23).

**Theorem 4.3.1.** *The optimal estimate  $\tilde{x}_m^{*,s}$  of  $x_m$  is given by the following expression*

$$\begin{aligned}
 \tilde{x}_m^{*,s} &= A_m \tilde{x}_{m-1}^{*,s} + \left[ (V_{m|m-1}^s)^{-1} + C_m^T (D_m D_m^T)^{-1} C_m \right]^{-1} C_m^T (D_m D_m^T)^{-1} \\
 &\times \left( y_m - C_m A_m \tilde{x}_{m-1}^{*,s} \right).
 \end{aligned} \tag{4.36}$$

Note that the above estimator can also be re-written by applying the PosDef Identity (see Appendix A) as

$$\tilde{x}_m^{*,s} = A_m \tilde{x}_{m-1}^{*,s} + V_{m|m-1}^s C_m^T \left[ C_m V_{m|m-1}^s C_m^T + (D_m D_m) \right]^{-1} \left( y_m - C_m A_m \tilde{x}_{m-1}^{*,s} \right). \tag{4.37}$$

This form is the same as the Kalman-Filter solution [1], [4]. The difference can be found in the covariance of the predictor error,  $V_{m|m-1}^s$ , which for this case it is a function of  $s$ .

Let

$$\begin{aligned}
 \hat{x}_{m|m} &\triangleq \lim_{s \rightarrow \infty} \hat{x}_{m|m}^s, \quad V_{m|m} \triangleq \lim_{s \rightarrow \infty} V_{m|m}^s, \quad \hat{x}_{m|m-1} \triangleq \lim_{s \rightarrow \infty} \hat{x}_{m|m-1}^s, \\
 V_{m|m-1} &\triangleq \lim_{s \rightarrow \infty} V_{m|m-1}^s, \quad \tilde{x}_m^* \triangleq \lim_{s \rightarrow \infty} \tilde{x}_m^{*,s}.
 \end{aligned}$$

Then, when  $s \rightarrow \infty$ ,

$$\hat{x}_{m|m} = V_{m|m}^s \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1})^{-1} \hat{x}_{m|m-1} \right] \quad (4.38)$$

$$V_{m|m} = \left( C_m^T (D_m D_m^T)^{-1} C_m + (V_{m|m-1})^{-1} \right)^{-1} \quad (4.39)$$

where

$$\hat{x}_{m|m-1} = A_m \hat{x}_{m-1|m-1} \quad (4.40)$$

$$V_{m|m-1} = B_m B_m^T + A_m V_{m-1|m-1} A_m^T \quad (4.41)$$

the estimator (4.37) converges to the classical Kalman-Filter solution, given by (1.11).

$$\tilde{x}_m^* = A_m \tilde{x}_{m-1}^* + V_{m|m-1} C_m^T \left[ C_m V_{m|m-1}^s C_m^T + (D_m D_m) \right] \left( y_m - C_m A_m \tilde{x}_{m-1}^* \right). \quad (4.42)$$

*Proof.* As was mentioned in Section 4.2.2 in order to find the optimal estimator  $\tilde{x}^{*,s}$  one needs to solve the recursive equation (4.13), then substitute into (4.21) and perform the minimization over  $\tilde{x}_m \in \mathcal{X}_{ad}$ . By applying the theory to this specific problem, then

$$\begin{aligned} \tilde{x}_m^{*,s} &= \arg \min_{\tilde{x}_m} s \log \int_{\mathbb{R}^{(m+1)d}} \alpha_m^s(1) dP_{y^m}(y) + sR \\ &= \arg \min_{\tilde{x}_m} \int_{\mathbb{R}^n} \bar{\alpha}_m^s(x) dx, \quad P_{y^m} - a.s. \\ &= \arg \min_{\tilde{x}_m} \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{\frac{n}{2}} |V_{m|m}^s|^{\frac{1}{2}}} \exp \left( - (x - \hat{x}_{m|m}^s)^T \frac{(V_{m|m}^s)^{-1}}{2} (x - \hat{x}_{m|m}^s) + \beta_{m|m}^s \right) dx \\ &= \arg \min_{\tilde{x}_m} \exp \left( \beta_{m|m}^s \right) \\ &= \arg \min_{\tilde{x}_m} \exp \left( (\tilde{x}_m - N_m)^T \frac{\Sigma_m^{-1}}{2} (\tilde{x}_m - N_m) + K(y_m, \hat{x}_{m-1|m-1}) \right) \end{aligned} \quad (4.43)$$

where

$$\begin{aligned} \Sigma_m^{-1} &= \frac{W_m}{s} V_{m|m}^s \frac{W_m}{s} + \frac{W_m}{s}, \\ N_m &= \Sigma_m \frac{W_m}{s} V_{m|m}^s \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1}^s)^{-1} \hat{x}_{m-1|m-1}^s \right]. \end{aligned}$$

Continuing from (4.43) and by performing the minimization

$$\tilde{x}_m^{*,s} = N_m$$

$$\begin{aligned}
 &= \left( \frac{W_m}{s} V_{m|m}^s \frac{W_m}{s} + \frac{W_m}{s} \right)^{-1} \frac{W_m}{s} V_{m|m}^s \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1}^s)^{-1} \hat{x}_{m|m-1}^s \right] \\
 &= \left( (V_{m|m}^s)^{-1} + \frac{W_m}{s} \right)^{-1} \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m}^s)^{-1} \hat{x}_{m|m-1}^s \right] \\
 &= \left( (V_{m|m-1}^s)^{-1} + C_m^T (D_m D_m^T)^{-1} C_m \right)^{-1} \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1}^s)^{-1} \hat{x}_{m|m-1}^s \right]
 \end{aligned} \tag{4.44}$$

$$= \hat{x}_{m|m-1}^s + \left( (V_{m|m-1}^s)^{-1} + C_m^T (D_m D_m^T)^{-1} C_m \right)^{-1} C_m^T (D_m D_m^T)^{-1} \left[ y_m - C_m \hat{x}_{m|m-1}^s \right]. \tag{4.45}$$

By using (4.28) and solving for  $\hat{x}_{m|m}^s$  (4.45) can be transformed into

$$\begin{aligned}
 \hat{x}_{m|m}^s &= V_{m|m}^s \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1}^s)^{-1} \hat{x}_{m|m-1}^s - \frac{W_m}{s} \tilde{x}_m^* \right] \\
 &= V_{m|m}^s \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1}^s)^{-1} \hat{x}_{m|m-1}^s \right] - V_{m|m}^s \frac{W_m}{s} \tilde{x}_m^*.
 \end{aligned} \tag{4.46}$$

Next, (4.44) can be expressed as

$$\left( (V_{m|m-1}^s)^{-1} + C_m^T (D_m D_m^T)^{-1} C_m \right) \tilde{x}_m^* = \left[ C_m^T (D_m D_m^T)^{-1} y_m + (V_{m|m-1}^s)^{-1} \hat{x}_{m|m-1}^s \right] \tag{4.47}$$

and when it is substituted into (4.46) it transforms it to

$$\begin{aligned}
 \hat{x}_{m|m}^s &= V_{m|m}^s \left( (V_{m|m}^s)^{-1} + C_m^T (D_m D_m^T)^{-1} C_m \right) \tilde{x}_m^{*,s} - V_{m|m}^s \frac{W_m}{s} \tilde{x}_m^{*,s} \\
 &= V_{m|m}^s \left( (V_{m|m-1}^s)^{-1} + C_m^T (D_m D_m^T)^{-1} C_m - \frac{W_m}{s} \right) \tilde{x}_m^{*,s} \\
 &= V_{m|m}^s (V_{m|m}^s)^{-1} \tilde{x}_m^{*,s} = \tilde{x}_m^{*,s}.
 \end{aligned} \tag{4.48}$$

Finally, (4.36) is derived by substituting (4.31) and (4.48) into (4.45).

□

## 4.4 Non-Coherent Estimation in Multipath

In this section the theory derived in Section 4.2 is applied to an attenuated sinusoid in a multipath environment which is subject to an additive Gaussian noise. Various

estimators are derived by solving the recursive equation for  $\{\bar{\alpha}_m^s(x) : m \in N_0\}$ . Further, as a special case, the estimation problem when there is no uncertainty is derived, by taking the limit as  $s \rightarrow \infty$ .

#### 4.4.1 Channel Model

A multipath version of the classical non-coherent model, given by the following equation, is considered in this section.

$$\begin{aligned} y(t_k) &= \sum_{i=1}^N A_i(t_k) r_i \cos(\omega_c(t_k - \tau_i(t_k)) + \theta_i) S(t_k - \tau_i(t_k)) + D(t_k) v(t_k) \\ &= \sum_{i=1}^N h_i(t_k, \theta_i, r_i) + D(t_k) v(t_k) \end{aligned} \quad (4.49)$$

where  $\omega_c$  is the carrier frequency,  $\{\tau_i(t_k)\}$  denotes the propagation delay,  $A_i(t_k)$  denotes a deterministic known signal envelope,  $\{r_i\}$ ,  $\{\theta_i\}$  are Random Variables denoting the attenuation and phase, respectively, of the signal received associated with  $i$ th path, and  $v(t_k) \sim N(0, 1)$ . Furthermore, the following function is defined  $h(N, \boldsymbol{\theta}, \mathbf{r}) \triangleq \sum_{i=1}^N h_i(t_k, \theta_i, r_i)$ , where  $\boldsymbol{\theta} \triangleq (\theta_1, \dots, \theta_N)'$  is the phase vector, and  $\mathbf{r} \triangleq (r_1, \dots, r_N)'$  is the attenuation vector.

The delays  $\{\tau_i(\cdot)\}_{i=1}^N$  are assumed to be fixed and known, while the phases  $\theta_i : \Omega \rightarrow [0, 2\pi]$  are independent and identically distributed RVs with *a priori* density  $\pi_{\theta_0}(\theta_i)$ ,  $\theta_i \in [0, 2\pi]$ , while the attenuations  $r_i : \Omega \rightarrow [0, \infty)$  are independent and identically distributed (iid) RVs with *a priori* density  $\pi_{r_0}(r_i)$ , for  $1 \leq i \leq N$ . In addition it is assumed  $\{r_i\}_{i=1}^N$  and  $\{\theta_i\}_{i=1}^N$  are independent, and also independent of the noise process  $\{v(t_k); k \in N_0\}$ . A special case of (4.49) in which  $N = 1$ ,  $\tau_1 = 0$  and  $r_1 = 1$  is the model used for non-coherent detection (see [4], [51]).

Further, it is assumed that the estimation error for  $\boldsymbol{\theta}$  and  $\mathbf{r}$  is additive given by

$$\ell(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}^m(t_k), \mathbf{r}, \tilde{\mathbf{r}}^m(t_k)) = \ell_1(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}^m(t_k)) + \ell_2(\mathbf{r}, \tilde{\mathbf{r}}^m(t_k))$$

where  $\ell_1(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}^m(t_k)) = \sum_{k=0}^m \sum_{i=1}^N \tilde{\ell}_{1,i}(\theta_i, \tilde{\theta}_i(t_k))$ ,  $\ell_2(\mathbf{r}, \tilde{\mathbf{r}}^m(t_k)) = \sum_{k=0}^m \sum_{i=1}^N \tilde{\ell}_{2,i}(r_i, \tilde{r}_i(t_k))$  and  $\tilde{\boldsymbol{\theta}}(t_k)$ ,  $\tilde{\mathbf{r}}(t_k)$  are the estimates of  $\boldsymbol{\theta}$ ,  $\mathbf{r}$ , respectively at time  $t_k$ ,  $k \in N_0$ . Note that  $\tilde{\ell}_{1,i}$ ,  $\tilde{\ell}_{2,i}$  are functions of estimation error.

The relation of model (4.49) with model (4.1) is the following. Clearly, sampling time  $k$  of previous section is now represented by  $t_k$ , and the state vector of (4.49) which

needs to be estimated is  $x = \begin{pmatrix} \boldsymbol{\theta} \\ \mathbf{r} \end{pmatrix} \in \mathfrak{R}^{2N}$ . However, since  $\begin{pmatrix} \boldsymbol{\theta} \\ \mathbf{r} \end{pmatrix}$  are random variables then  $x_{k+1} = x_k, \forall k \in N_0$ .

Next, the theory of Section 4.2 is applied to the above problem.

#### 4.4.2 Minimax Estimation of Phase and Envelope

The density (4.16) is specialized to model (4.49). This is done by replacing the integrand  $\Psi_w(B_m^{-1}(x_m - f(t, z)))$ ,  $x \in \mathfrak{R}^{2N}$ ,  $z \in \mathfrak{R}^{2N}$  in (4.16) with the delta measure  $\delta(\mathbf{r} - z^1) \times \delta(\boldsymbol{\theta} - z^2) = \prod_{i=1}^N \delta(r_i - z_i^1) \times \delta(\theta_i - z_i^2)$ , where  $x = \begin{pmatrix} \boldsymbol{\theta} \\ \mathbf{r} \end{pmatrix}$ ,  $z = \begin{pmatrix} z^1 \\ z^2 \end{pmatrix} \in \mathfrak{R}^{2N}$ ,  $f(t, z) = \begin{pmatrix} z^1 \\ z^2 \end{pmatrix}$  to get

$$\begin{aligned} \bar{\alpha}_m^s(\boldsymbol{\theta}, \mathbf{r}) &= \pi_{\boldsymbol{\theta}_0}(\boldsymbol{\theta}) \pi_{\mathbf{r}_0}(\mathbf{r}) \prod_{k=0}^m \left[ \frac{\Xi_v(D^{-1}(t_k)(y(t_k) - \sum_{i=1}^N h_i(t_k, \theta_i, r_i)))}{|D(t_k)| \Xi_v(y(t_k))} \right. \\ &\quad \left. \times \exp \left( \frac{1}{s} \sum_{i=1}^N \tilde{\ell}_{1,i}(\theta_i, \tilde{\theta}_i(t_k)) + \frac{1}{s} \sum_{i=1}^N \tilde{\ell}_{2,i}(r_i, \tilde{r}_i(t_k)) \right) \right] \end{aligned} \quad (4.50)$$

where  $\pi_{\boldsymbol{\theta}_0}(\boldsymbol{\theta}) = \prod_{i=1}^N \pi_{\theta_0}(\theta_i)$  is the *a priori* joint density of  $(\theta_1, \dots, \theta_N)$ , and  $\pi_{\mathbf{r}_0}(\mathbf{r}) = \prod_{i=1}^N \pi_{r_0}(r_i)$  is the *a priori* joint density of  $(r_1, \dots, r_N)$ .

Furthermore, the error functions are assumed to have the form  $\tilde{\ell}_{1,i}(\theta_i, \tilde{\theta}_i(t_k)) = \frac{(\theta_i - \tilde{\theta}_i(t_k))^2 R_i(t_k)}{2}$ , and  $\tilde{\ell}_{2,i}(r_i, \tilde{r}_i(t_k)) = \frac{(r_i - \tilde{r}_i(t_k))^2 Q_i(t_k)}{2}$ , where  $R_i(\cdot)$ ,  $Q_i(\cdot)$  are weighting coefficients on the specific errors, which are positive for each  $t_k$ .

The following quantities are needed when presenting the robust estimator.

**Definition 4.4.1.** For each  $1 \leq i, j \leq N, i \neq j$ , define the following quantities.

$$\begin{aligned} V_c^i(y^m) &\triangleq \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \cos(\omega_c(t_k - \tau_i(t_k))) S(t_k - \tau_i(t_k)) y(t_k). \\ V_s^i(y^m) &\triangleq \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \sin(\omega_c(t_k - \tau_i(t_k))) S(t_k - \tau_i(t_k)) y(t_k). \\ V_i(y^m) &= \sqrt{V_c^i(y^m)^2 + V_s^i(y^m)^2}, \quad \gamma_i(y^m) = -\tan^{-1} \left( \frac{V_s^i(y^m)}{V_c^i(y^m)} \right). \\ W_c^{i,m} &\triangleq -\frac{1}{4} \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \cos(2\omega_c(t_k - \tau_i(t_k))) S^2(t_k - \tau_i(t_k)). \\ W_s^{i,m} &\triangleq -\frac{1}{4} \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \sin(2\omega_c(t_k - \tau_i(t_k))) S^2(t_k - \tau_i(t_k)). \\ W_i^m &= \sqrt{(W_c^{i,m})^2 + (W_s^{i,m})^2}, \quad \beta_i^m = -\tan^{-1} \left( \frac{W_s^{i,m}}{W_c^{i,m}} \right). \end{aligned}$$



$$\begin{aligned}
 U_c^{ij,m} &\triangleq -\frac{1}{2} \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \cos(\omega_c(2t_k - \tau_i(t_k) - \tau_j(t_k))) S(t_k - \tau_i(t_k)) S(t_k - \tau_j(t_k)). \\
 U_s^{ij,m} &\triangleq -\frac{1}{2} \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \sin(\omega_c(2t_k - \tau_i(t_k) - \tau_j(t_k))) S(t_k - \tau_i(t_k)) S(t_k - \tau_j(t_k)). \\
 U_{ij}^m &= \sqrt{(U_c^{ij,m})^2 + (U_s^{ij,m})^2}, \quad \phi_{ij}^m = -\tan^{-1} \left( \frac{U_s^{ij,m}}{U_c^{ij,m}} \right). \\
 T_c^{ij,m} &\triangleq -\frac{1}{2} \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \cos(\omega_c(\tau_j(t_k) - \tau_i(t_k))) S(t_k - \tau_i(t_k)) S(t_k - \tau_j(t_k)). \\
 T_s^{ij,m} &\triangleq -\frac{1}{2} \sum_{k=0}^m A_i(t_k) D^{-2}(t_k) \sin(\omega_c(\tau_j(t_k) - \tau_i(t_k))) S(t_k - \tau_i(t_k)) S(t_k - \tau_j(t_k)). \\
 T_{ij}^m &= \sqrt{(T_c^{ij,m})^2 + (T_s^{ij,m})^2}, \quad \psi_{ij}^m = -\tan^{-1} \left( \frac{T_s^{ij,m}}{T_c^{ij,m}} \right).
 \end{aligned}$$

From Definition 4.4.1 and (4.50) the unnormalized conditional density is given by

$$\begin{aligned}
 \bar{\alpha}_m^s(\boldsymbol{\theta}, \mathbf{r}) &= \pi_{\boldsymbol{\theta}_0}(\boldsymbol{\theta}) \pi_{\mathbf{r}_0}(\mathbf{r}) \exp \left( -\frac{1}{4} \sum_{i=1}^N \sum_{k=0}^m r_i^2 A_i^2(t_k) D^{-2}(t_k) S^2(t_k - \tau_i(t_k)) \right) \\
 &\times \exp \left( \sum_{i=1}^N r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) \exp \left( \sum_{i=1}^N r_i^2 W_i^m \cos(2\theta_i - \beta_i^m) \right) \\
 &\times \exp \left( \sum_{i=1}^N \sum_{j=i+1}^N r_i r_j U_{ij}^m \cos(\theta_i + \theta_j - \phi_{ij}^m) \right) \\
 &\times \exp \left( \sum_{i=1}^N \sum_{j=i+1}^N r_i r_j T_{ij}^m \cos(\theta_i - \theta_j - \psi_{ij}^m) \right) \\
 &\times \exp \left( \sum_{i=1}^N \sum_{k=0}^m \frac{(\theta_i - \tilde{\theta}_i(t_k))^2 R_i(t_k)}{2s} + \sum_{i=1}^N \sum_{k=0}^m \frac{(r_i - \tilde{r}_i(t_k))^2 Q_i(t_k)}{2s} \right) \\
 &\times \exp \left( \frac{1}{2} \sum_{k=0}^m y^2(t_k) [1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)| \right). \tag{4.51}
 \end{aligned}$$

Note that (4.51) should be substituted into (4.21) and then minimized over the class of estimators. Since the exponential term is a quadratic function of  $(\boldsymbol{\theta}, \mathbf{r})$  and  $(\tilde{\boldsymbol{\theta}}^m, \tilde{\mathbf{r}}^m)$ , the minimization in (4.51) over  $\{\tilde{\boldsymbol{\theta}}^m, \tilde{\mathbf{r}}^m\}$  can be found explicitly provided  $\pi_{\boldsymbol{\theta}_0}(\cdot)$ ,  $\pi_{\mathbf{r}_0}(\cdot)$  are Gaussian. However, for this problem this is not the case. Nevertheless, closed form expressions for the optimal estimators when  $\boldsymbol{\theta}$  is uniformly distributed over  $[0, 2\pi]$  and  $\mathbf{r}$  is arbitrary distributed are going to be derived.

**Remark 4.4.2.** *Below realistic scenarios that lead to a simplified version of (4.51) are discussed.*

1. *Minimax Estimation Subject to Resolvable Paths.*

*Supposed the received paths are resolvable, where path resolvability for wide-band signals is defined as having inter-path delays which are larger than the*

reciprocal of the bandwidth of the transmitted signals. Under the resolvability assumption, expressions containing cross terms in Definition 4.4.1 are negligible. This implies that the following terms are zero.

$$U_c^{ij,m} = U_s^{ij,m} = 0, \quad T_c^{ij,m} = T_s^{ij,m} = 0, \quad 1 \leq i, j \leq N, i \neq j.$$

Taking this into consideration, expression (4.51) reduces to

$$\begin{aligned} \bar{\alpha}_m^s(\boldsymbol{\theta}, \mathbf{r}) &= \pi_{\boldsymbol{\theta}_0}(\boldsymbol{\theta})\pi_{\mathbf{r}_0}(\mathbf{r}) \exp\left(-\frac{1}{4} \sum_{i=1}^N \sum_{k=0}^m r_i^2 A_i^2(t_k) D^{-2}(t_k) S^2(t_k - \tau_i(t_k))\right) \\ &\times \exp\left(\sum_{i=1}^N r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m))\right) \\ &\times \exp\left(\sum_{i=1}^N r_i^2 W_i^m \cos(2\theta_i - \beta_i^m)\right) \\ &\times \exp\left(\sum_{i=1}^N \sum_{k=0}^m \frac{(\theta_i - \tilde{\theta}_i(t_k))^2 R_i(t_k)}{2s} + \sum_{i=1}^N \sum_{k=0}^m \frac{(r_i - \tilde{r}_i(t_k))^2 Q_i(t_k)}{2s}\right) \\ &\exp\left(\frac{1}{2} \sum_{k=0}^m y^2(t_k) [1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)|\right). \end{aligned} \quad (4.52)$$

## 2. Minimax Estimation Subject to Negligible Double Frequency Terms.

Here a ‘‘Narrow-band’’ assumption is employed, which implies  $2\omega_c$  terms (‘‘double-frequency’’ terms) are negligible since the receiver will remove them in the process of reconstructing the transmitted signal.

Taking this into consideration, expression (4.51) reduces to

$$\begin{aligned} \bar{\alpha}_m^s(\boldsymbol{\theta}, \mathbf{r}) &= \pi_{\boldsymbol{\theta}_0}(\boldsymbol{\theta})\pi_{\mathbf{r}_0}(\mathbf{r}) \exp\left(-\frac{1}{4} \sum_{i=1}^N \sum_{k=0}^m r_i^2 A_i^2(t_k) D^{-2}(t_k) S^2(t_k - \tau_i(t_k))\right) \\ &\times \exp\left(\sum_{i=1}^N r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m))\right) \\ &\times \exp\left(\sum_{i=1}^N \sum_{j=i+1}^N r_i r_j U_{ij}^m \cos(\theta_i + \theta_j - \phi_{ij}^m)\right) \\ &\times \exp\left(\sum_{i=1}^N \sum_{j=i+1}^N r_i r_j T_{ij}^m \cos(\theta_i - \theta_j - \psi_{ij}^m)\right) \\ &\times \exp\left(\sum_{i=1}^N \sum_{k=0}^m \frac{(\theta_i - \tilde{\theta}_i(t_k))^2 R_i(t_k)}{2s} + \sum_{i=1}^N \sum_{k=0}^m \frac{(r_i - \tilde{r}_i(t_k))^2 Q_i(t_k)}{2s}\right) \\ &\times \exp\left(\frac{1}{2} \sum_{k=0}^m y^2(t_k) [1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)|\right). \end{aligned} \quad (4.53)$$

If cases (1), (2) above are combined together the density  $\bar{\alpha}_m^s(\boldsymbol{\theta}, \mathbf{r})$  is given by (4.53) with cross terms  $U_c^{ij,m}, U_s^{ij,m}, T_c^{ij,m}, T_s^{ij,m}$  being zero. This is treated in the following paragraphs.

Supposed the received paths are resolvable, and double frequency terms  $2\omega_c$  are negligible. Then (4.51) reduces to

$$\bar{\alpha}_m^s(\boldsymbol{\theta}, \mathbf{r}) = \left\{ \prod_{i=1}^N \bar{\alpha}_m^{s,i}(\theta_i, r_i) \right\} \times \exp \left( \frac{1}{2} \sum_{k=0}^m y^2(t_k) [1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)| \right) \quad (4.54)$$

where  $\bar{\alpha}_m^{s,i}(\theta_i, r_i)$ ,  $1 \leq i \leq N$  is given by

$$\begin{aligned} \bar{\alpha}_m^{s,i}(\theta_i, r_i) &= \pi_{\theta_0}(\theta_i) \pi_{r_0}(r_i) \exp \left( -\frac{1}{4} \sum_{k=0}^m r_i A_i^2(t_k) D^{-2}(t_k) S^2(t_k - \tau_i(t_k)) \right) \\ &\times \exp \left( r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) \\ &\times \exp \left( \frac{(\theta_i - \tilde{\theta}_i(t_m))^2 R_i(t_m)}{2s} + \sum_{k=0}^{m-1} \frac{(\theta_i - \tilde{\theta}_i(t_k))^2 R_i(t_k)}{2s} \right) \\ &\times \exp \left( \frac{(r_i - \tilde{r}_i(t_m))^2 Q_i(t_m)}{2s} + \sum_{k=0}^{m-1} \frac{(r_i - \tilde{r}_i(t_k))^2 Q_i(t_k)}{2s} \right). \end{aligned} \quad (4.55)$$

Note that path resolvability implies that each path component  $(\theta_i, r_i)$  can be estimated independently of the rest.

**Theorem 4.4.3.** *Assuming negligible double frequency terms, and resolvable paths, then the optimal robust estimates  $\{\tilde{\theta}_i^{*,s}(t_m)\}_{i=1}^N$  and  $\{\tilde{r}_i^{*,s}(t_m)\}_{i=1}^N$  are given by the following expressions*

$$\begin{aligned} \tilde{\theta}_i^{*,s}(t_m) &= \\ &\frac{\int_0^\infty \int_0^{2\pi} \theta_i \pi_{r_0}(r_i) \exp(-r_i^2 K_i^m) \exp \left( r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) \exp(\zeta_i^m + \xi_i^m) d\theta_i dr_i}{\int_0^\infty \int_0^{2\pi} \pi_{r_0}(r_i) \exp(-r_i^2 K_i^m) \exp \left( r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) \exp(\zeta_i^m + \xi_i^m) d\theta_i dr_i} \end{aligned} \quad (4.56)$$

$$\begin{aligned} \tilde{r}_i^{*,s}(t_m) &= \\ &\frac{\int_0^\infty \int_0^{2\pi} r_i \pi_{r_0}(r_i) \exp(-r_i^2 K_i^m) \exp \left( r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) \exp(\zeta_i^m + \xi_i^m) d\theta_i dr_i}{\int_0^\infty \int_0^{2\pi} \pi_{r_0}(r_i) \exp(-r_i^2 K_i^m) \exp \left( r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) \exp(\zeta_i^m + \xi_i^m) d\theta_i dr_i} \end{aligned} \quad (4.57)$$

where

$$\zeta_i^m = \frac{(\theta_i - \tilde{\theta}_i^{*,s}(t_m))^2 R_i(t_m)}{2s} + \sum_{k=0}^{m-1} \frac{(\theta_i - \tilde{\theta}_i^{*,s}(t_k))^2 R_i(t_k)}{2s}, \quad (4.58)$$

$$\xi_i^m = \frac{(r_i - \tilde{r}_i^{*,s}(t_m))^2 Q_i(t_m)}{2s} + \sum_{k=0}^{m-1} \frac{(r_i - \tilde{r}_i^{*,s}(t_k))^2 Q_i(t_k)}{2s}, \quad (4.59)$$

$$K_i^m = \frac{1}{4} \sum_{k=0}^m A_i^2(t_k) D^{-2}(t_k) S^2(t_k - \tau_i(t_k)). \quad (4.60)$$

*Proof.* Substituting (4.54) into (4.21) and using the fact that  $P_{y^m}(y)$  is a positive density function, which is independent of  $\tilde{\theta}_i$  and  $\tilde{r}_i$ , then  $\tilde{\theta}_i^{*,s}$  and  $\tilde{r}_i^{*,s}$  are found from the expression

$$\begin{aligned} \{(\tilde{\theta}_i^{*,s}(t_m), \tilde{r}_i^{*,s}(t_m))\}_{i=1}^N &= \arg \min_{\{(\tilde{\theta}_i(t_m), \tilde{r}_i(t_m))\}_{i=1}^N} \log \prod_{i=1}^N \left[ \int_0^\infty \int_0^{2\pi} \bar{\alpha}_m^{s,i}(\theta_i, r_i) d\theta_i dr_i \right] \\ &= \arg \min_{\{(\tilde{\theta}_i(t_m), \tilde{r}_i(t_m))\}_{i=1}^N} \sum_{i=1}^N \left[ \log \int_0^\infty \int_0^{2\pi} \bar{\alpha}_m^{s,i}(\theta_i, r_i) d\theta_i dr_i \right]. \end{aligned} \quad (4.61)$$

Hence

$$(\tilde{\theta}_i^{*,s}(t_m), \tilde{r}_i^{*,s}(t_m)) = \arg \min_{(\tilde{\theta}_i(t_m), \tilde{r}_i(t_m))} \log \int_0^\infty \int_0^{2\pi} \bar{\alpha}_m^{s,i}(\theta_i, r_i) d\theta_i dr_i, \quad 1 \leq i \leq N. \quad (4.62)$$

Differentiating the pay-off function (4.62) with respect to  $\tilde{\theta}_i(t_m)$  and  $\tilde{r}_i(t_m)$ , evaluating at  $(\tilde{\theta}_i(t_m), \tilde{r}_i(t_m)) = (\tilde{\theta}_i^{*,s}(t_m), \tilde{r}_i^{*,s}(t_m))$ , setting it to zero and then solving for  $(\tilde{\theta}_i^{*,s}(t_m), \tilde{r}_i^{*,s}(t_m))$  yields expressions (4.56), (4.57).

Notice that the minimax estimators of the phase  $\tilde{\theta}_i^{*,s}(t_m)$  and envelope  $\tilde{r}_i^{*,s}(t_m)$  are computed from (4.56) and (4.57). Clearly, the estimators derived in (4.56) and (4.57) are given implicitly because  $\tilde{\theta}_i^{*,s}(t_m)$  and  $\tilde{r}_i^{*,s}(t_m)$  appear on both left and right side of these expressions. In calculating these quantities, one should start with an initial guess  $\tilde{\theta}_i^{*,s} = \tilde{\theta}_i^{1,s}$  and  $\tilde{r}_i^{*,s} = \tilde{r}_i^{1,s}$ , put them into the right side of (4.56) and (4.57) compute  $\tilde{\theta}_i^{2,s}$ ,  $\tilde{r}_i^{2,s}$  and repeat this step until the estimators converge.  $\square$

### 4.4.3 Classical Estimation of Phase and Envelope

Here the problem of jointly estimating the channel parameters, namely phase,  $\boldsymbol{\theta}$  and envelope,  $\mathbf{r}$  for an attenuated sinusoid signal in a multipath environment which is subjected to additive Gaussian noise, when there is no uncertainty, is considered. The case of no uncertainty corresponds to having an uncertainty radius  $R = 0$ . From Theorem 4.2.2, this is equivalent to taking the limit, as  $s \rightarrow \infty$  [43]. Thus, the limit  $\bar{\alpha}_m^\infty(\boldsymbol{\theta}, \mathbf{r}) \triangleq \lim_{s \rightarrow \infty} \bar{\alpha}_m^s(\boldsymbol{\theta}, \mathbf{r})$  corresponds to the unnormalized conditional density of  $(\boldsymbol{\theta}, \mathbf{r})$  given  $\mathcal{F}_m^y$ . The mean-square error estimation of  $(\boldsymbol{\theta}, \mathbf{r})$  is then computed via

$$\tilde{\boldsymbol{\theta}}^*(t_m) = E[\boldsymbol{\theta} | \mathcal{F}_m^y] = \frac{\int \boldsymbol{\theta} \bar{\alpha}_m^\infty(\boldsymbol{\theta}, \mathbf{r}) d\boldsymbol{\theta} d\mathbf{r}}{\int \bar{\alpha}_m^\infty(\boldsymbol{\theta}, \mathbf{r}) d\boldsymbol{\theta} d\mathbf{r}},$$

and similarly for  $\tilde{\mathbf{r}}^*(t_m)$ .

From (4.54), since it is assumed negligible double frequency terms and resolvable paths, and by reducing the uncertainty to zero (by letting  $s \rightarrow \infty$ ), the unnormalized conditional density of  $(\boldsymbol{\theta}, \mathbf{r})$  given  $\mathcal{F}_m^y$  is given by

$$\begin{aligned} \bar{\alpha}_m^\infty(\boldsymbol{\theta}, \mathbf{r}) &= \prod_{i=1}^N \left[ \pi_{\theta_0}(\theta_i) \pi_{r_0}(r_i) \exp(-r_i^2 K_i^m) \exp\left(r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m))\right) \right] \\ &\times \exp\left(\frac{1}{2} \sum_{k=0}^m y^2(t_k) [1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)|\right) \end{aligned} \quad (4.63)$$

where  $K_i^m = \frac{1}{4} \sum_{k=0}^m A_i^2(t_k) D^{-2}(t_k) S^2(t_k - \tau_i(t_k))$ .

Note that path resolvability implies that each path component  $(\theta_i, r_i)$  can be estimated independently of the rest. Below various optimal estimator are derived.

**Theorem 4.4.4.** *Assuming the phases  $\theta_i$  are iid RVs with a priori density  $\pi_{\theta_0}(\theta_i) = \frac{1}{2\pi}$ ,  $\theta_i \in [0, 2\pi]$  and the attenuations  $r_i$  are iid RVs with a priori density  $\pi_{r_0}(r_i) = \frac{r_i}{\sigma^2} \exp -\frac{r_i^2}{2\sigma^2}$ ,  $r_i \in [0, \infty)$  (Rayleigh distributed), then the following estimators are obtained.*

(a) *The incomplete data likelihood ratio defined by  $\hat{\Lambda}(t_m) = \bar{E}[\Lambda(t_m) | \mathcal{F}_m^y]$  is given by*

$$\begin{aligned} \hat{\Lambda}(t_m) &= \alpha_m(1) = \int_{[0, \infty)^n} \int_{[0, 2\pi]^n} \bar{\alpha}_m(\boldsymbol{\theta}, \mathbf{r}) d\boldsymbol{\theta} d\mathbf{r} \\ &= \prod_{i=1}^N \left[ \frac{1}{1 + 2\sigma^2 K_i^m} \exp\left(\frac{V_i^2(y^m) \sigma^2}{2 + 4\sigma^2 K_i^m}\right) \right] \end{aligned} \quad (4.64)$$

$$\times \exp \left( \frac{1}{2} \sum_{k=0}^m y^2(t_k) [1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)| \right). \quad (4.65)$$

(b) The normalized conditional density of  $(\boldsymbol{\theta}, \mathbf{r})$  given  $\mathcal{F}_m^{\mathcal{Y}}$ , i.e.,  $p_N(t_m, \boldsymbol{\theta}, \mathbf{r} | \mathcal{F}_m^{\mathcal{Y}})$ , is given by

$$\begin{aligned} p_N(t_m, \boldsymbol{\theta}, \mathbf{r} | \mathcal{F}_m^{\mathcal{Y}}) &= \frac{\bar{\alpha}_m(\boldsymbol{\theta}, \mathbf{r})}{\int \int \bar{\alpha}_m(\boldsymbol{\theta}, \mathbf{r}) d\boldsymbol{\theta} d\mathbf{r}} = \prod_{i=1}^N p_N(t_m, \theta_i, r_i | \mathcal{F}_m^{\mathcal{Y}}) \quad (4.66) \\ &= \prod_{i=1}^N \left[ \frac{r_i(1 + 2\sigma^2 K_i^m)}{2\pi\sigma^2} \exp \left( -\frac{r_i^2(1 + 2\sigma^2 K_i^m)}{2\sigma^2} \right) \right. \\ &\quad \times \exp \left( r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) \\ &\quad \left. \times \exp \left( -\frac{V_i^2(y^m)\sigma^2}{2 + 4\sigma^2 K_i^m} \right) \right]. \quad (4.67) \end{aligned}$$

(c) The conditional least-square estimate of the noiseless received signal  $h_i(t_m, \theta_i, r_i)$  given  $\mathcal{F}_m^{\mathcal{Y}}$ , defined by  $\hat{h}_i(t_m, \theta_i, r_i) \triangleq E[h_i(t_m, \theta_i, r_i) | \mathcal{F}_m^{\mathcal{Y}}]$  is given by

$$\begin{aligned} \hat{h}_i(t_m, \theta_i, r_i) &= \int_0^\infty \int_0^{2\pi} h_i(t_m, \theta_i, r_i) p_N(t_m, \theta_i, r_i | \mathcal{F}_m^{\mathcal{Y}}) d\theta_i dr_i \quad (4.68) \\ &= S(t_m - \tau_i(t_m)) \cos(\omega_c(t_m - \tau_i(t_m)) + \gamma_i(y^m)) \frac{V_i(y^m)\sigma^2}{1 + 2\sigma^2 K_i^m}. \quad (4.69) \end{aligned}$$

(d) The minimum least-square estimator of  $\theta_i$  given  $\mathcal{F}_m^{\mathcal{Y}}$  is given by

$$\begin{aligned} \tilde{\theta}_i^{*,\infty}(t_m) &= E[\theta_i | \mathcal{F}_m^{\mathcal{Y}}] = \int_0^\infty \int_0^{2\pi} \theta_i p_N(t_m, \theta_i, r_i | \mathcal{F}_m^{\mathcal{Y}}) d\theta_i dr_i \quad (4.70) \\ &= \frac{(1 + 2\sigma^2 K_i^m)}{2\pi\sigma^2} \exp \left( -\frac{V_i^2(y^m)\sigma^2}{2 + 4\sigma^2 K_i^m} \right) \\ &\quad \times \int_0^\infty \int_0^{2\pi} \theta_i r_i \exp \left( -\frac{r_i^2(1 + 2\sigma^2 K_i^m)}{2\sigma^2} \right) \\ &\quad \times \exp \left( r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m)) \right) d\theta_i dr_i. \quad (4.71) \end{aligned}$$

(e) The minimum least-square estimator of  $r_i$  given  $\mathcal{F}_m^{\mathcal{Y}}$  is given by

$$\begin{aligned} \tilde{r}_i^{*,\infty}(t_m) &= E[r_i | \mathcal{F}_m^{\mathcal{Y}}] = \int_0^\infty \int_0^{2\pi} r_i p_N(t_m, \theta_i, r_i | \mathcal{F}_m^{\mathcal{Y}}) dr_i d\theta_i \quad (4.72) \\ &= (\sqrt{\pi/2}) \left( \sqrt{\frac{\sigma^2}{1 + 2\sigma^2 K_i^m}} \right) \exp \left( -\frac{V_i^2(y^m)\sigma^2}{2 + 4\sigma^2 K_i^m} \right) {}_1F_1 \left( \frac{3}{2}, 1; \frac{V_i^2(y^m)\sigma^2}{2 + 4\sigma^2 K_i^m} \right) \quad (4.73) \end{aligned}$$

where  ${}_1F_1(\alpha, \beta; x)$  is the confluent hypergeometric function [52].

*Proof.* See Appendix D. □

**Remark 4.4.5.** The results of Theorem 4.4.4 can be used in decision problems and in nonlinear estimation problems.

The detection problem associated with (4.49) is described by the following binary hypothesis problem.

$$\begin{aligned} H_1 : y(t_k) &= \sum_{i=1}^N \left[ A_i(t_k) r_i \cos(\omega_c(t_k - \tau_i(t_k)) + \theta_i) S(t_k - \tau_i(t_k)) \right] + D(t_k)v(t_k) \\ H_2 : y(t_k) &= D(t_k)v(t_k). \end{aligned}$$

The incomplete data likelihood ratio  $\hat{\Lambda}(t_m)$  can be used in minimum risk Bayes' decision applications, where there will be one distribution  $\alpha_m(x)$  for each hypothesis. That is, in binary hypothesis testing the decision rule takes the form

$$\frac{\int_{\mathbb{R}^n} d\alpha_m^1(x)}{\int_{\mathbb{R}^n} d\alpha_m^2(x)} = \frac{\hat{\Lambda}^1(t_m)}{\hat{\Lambda}^2(t_m)} \underset{H_2}{\overset{H_1}{\geq}} \gamma, \quad (4.74)$$

where  $\gamma$  is the threshold. Clearly,  $\alpha_m(x)$  is a sufficient statistic for estimation and decision problems, since both the least-square estimation and decision rule are constructed from this quantity.

The normalized conditional density  $p_N(t_m, x | \mathcal{F}_m^y)$  can be used to compute various least-square estimates as was shown in (4.71) and (4.73). Hence, the general expression for the least-square estimate can be written as

$$E[\Phi(x_m) | \mathcal{F}_m^y] = \int_{\mathbb{R}^n} \Phi(x_m) p_N(t_m, x | \mathcal{F}_m^y) dx. \quad (4.75)$$

Notice that from (4.69) the following least-squares estimate can also be derived

$$E \left[ r_i \cos(\omega_c(t_k - \tau_i(t_k)) + \theta_i) | \mathcal{F}_m^y \right] = \cos(\omega_c(t_m - \tau_i(t_m)) + \gamma_i(y^m)) \frac{V_i(y^m) \sigma^2}{1 + 2\sigma^2 K_i^m}. \quad (4.76)$$

Both estimators, (4.69) and (4.76), can be used in RAKE receiver applications. In a RAKE receiver, one RAKE finger is assigned to each multipath, thus maximizing the amount of received signal energy. Each of these different paths are combined to form

a composite signal which has substantially better characteristics for the purpose of demodulation, when compared to a single path. In previous work [52] in order to combine the different paths meaningfully, the RAKE receiver needs the knowledge of the channel parameters. However, using (4.69), (4.76) a RAKE receiver can be constructed without knowledge of the channel parameters.

#### 4.4.4 Numerical Results and Discussion

Here the performance of the estimators derived in Sections 4.4.2 and 4.4.3 is evaluated. Their performance is evaluated through the Mean Square Error function of each estimator for each arriving path as a function of time for  $N = 100$  realizations. The general expression of the MSE is given by

$$MSE_{Z_i} = \frac{1}{N} \sum_{j=1}^N |Z_i - \tilde{Z}_{i,j}^*|^2, \quad (4.77)$$

where  $N$  is the number of realizations,  $Z_i$  is the real value of the parameter that would be estimated for path  $i$  and  $\tilde{Z}_{i,j}^*$  is the estimated value of that parameter for path  $i$ .

The experiments presented here are for evaluation purposes and do not represent real life scenarios. A set of observations  $y$  is created for each case based on the non-coherent model given by (4.49). This model represents the baseband signal and all the various parameters, except the phase  $\theta$  and attenuation  $r$ , are assumed to be known. Therefore, a relatively small carrier frequency,  $f_c = 1\text{KHz}$ , will be used through this section together with a transmitted signal  $S(t_k) = 1$  and a Gaussian noise  $v(t_k) \sim N(0, \sigma_n^2)$ . Finally, the Signal to Noise Ratio is defined as  $SNR = \frac{P_s}{\sigma_n^2}$ , where  $P_s$  is the power of the transmitted signal.

For all the experiments presented below Wolfram Mathematica is being used.

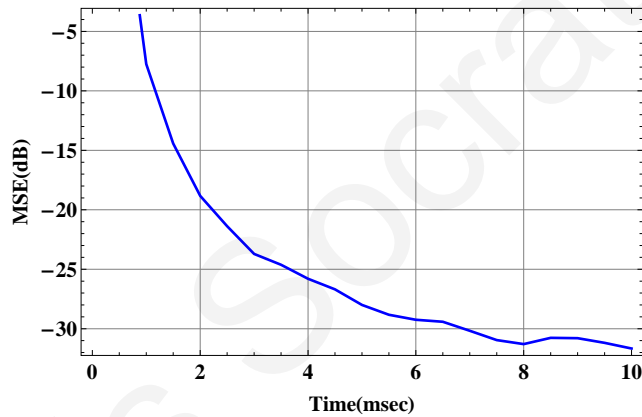
##### *Minimax Non-Coherent Model*

Firstly, the robust phase and envelope estimators derived in Section 4.4.2 are evaluated, that is the derived phase estimator (4.56) and attenuation estimator (4.57).

A single path version of the non-coherent model given by (4.49) is used and a set of observations  $y$  is created with reference to this model. The true phase  $\theta$  and



attenuation  $r$  parameters are random variables, rayleigh distributed, with  $\hat{a}$  priori densities  $\pi_\theta(\theta) = \frac{\theta}{\sigma_\theta^2} \exp\left(-\frac{\theta^2}{2\sigma_\theta^2}\right)$ ,  $\theta \in [0, 2\pi]$  and  $\pi_r(r) = \frac{r}{\sigma_r^2} \exp\left(-\frac{r^2}{2\sigma_r^2}\right)$ ,  $r \in [0, \infty)$ , respectively. The attenuation parameters of the above rayleigh distributions are taken as  $\sigma_\theta = 1.5$  and  $\sigma_r = \frac{1}{\sqrt{2}}$ , respectively. The true distributions of the phase and the attenuation are unknown to the observer and the estimators use the nominal distributions which are assumed to be a uniform distribution for the phase, with  $\hat{a}$  priori density  $\pi_{\theta_0}(\theta_i) = \frac{1}{2\pi}$ ,  $\theta_i \in [0, 2\pi]$ , and a rayleigh distribution for the attenuation, with an attenuation parameter of  $\sigma_r = \frac{5}{\sqrt{2}}$ . As was mentioned above, the performance of the derived estimators will be evaluated using the Mean Square Error (4.77) of each estimator over a period of time of  $10\text{msec}$ , for  $SNR = 20\text{dB}$ , and  $N = 100$  realizations.

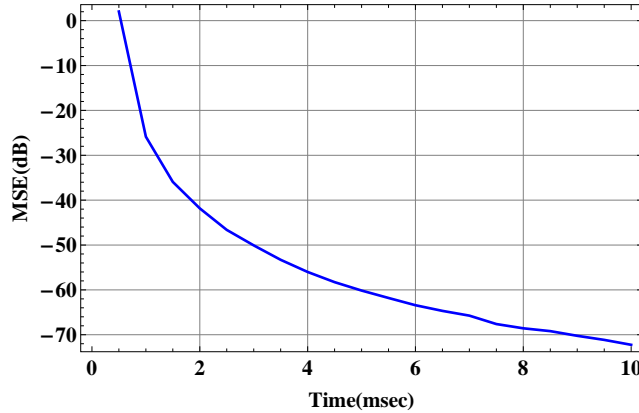


**Figure 4.1:** MSE of robust minimax phase estimator for a reference  $SNR=20$  dB

Starting with the robust phase estimation, Figure 4.1 displays the performance of the phase estimator (4.56) over time for a reference  $SNR = 20\text{dB}$ . Clearly the MSE decreases with time, as time increases and this shows that the estimated phase converges in mean square sense to the real phase.

Next, the performance of the robust envelope estimator (4.57) is evaluated. Figure 4.2 displays the performance of the attenuation estimator over time for a reference  $SNR = 20\text{dB}$ . The observations are similar with the ones experienced for the robust phase MSE and show that the MSE decreases as time increases and that the estimated attenuation converges in mean square sense to its real value.

Therefore, even though the nominal distributions that were chosen for both the phase and the attenuation were different from the true ones the estimated parameters converges in mean square sense to their real values.



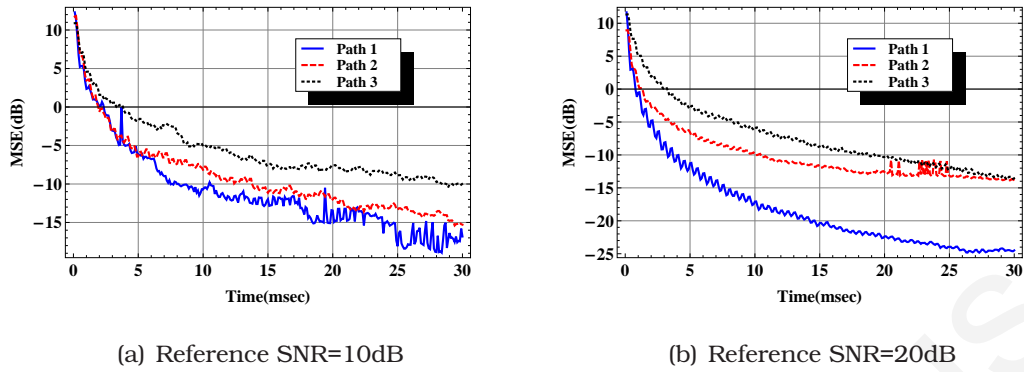
**Figure 4.2:** MSE of robust minimax attenuation estimator for a reference SNR=20 dB

#### *Classical Non-Coherent Model*

Next the phase, envelope and noiseless received signal estimators derived in Section 4.4.3 are evaluated. A multipath version of the classical non-coherent model, given by (4.49) with three resolvable paths, is considered here. It is assumed that the phases  $\theta_i$  are iid random variables with  $\hat{a}$  priori density  $\pi_{\theta_0}(\theta_i) = \frac{1}{2\pi}$ ,  $\theta_i \in [0, 2\pi]$  and the attenuations  $r_i$  are iid random variables with  $\hat{a}$  priori density  $\pi_{r_0}(r_i) = \frac{r_i}{\sigma_{r_i}^2} \exp\left(-\frac{r_i^2}{2\sigma_{r_i}^2}\right)$ ,  $r_i \in [0, \infty)$  (Rayleigh distributed). It is also assumed that the attenuation parameter  $\sigma_{r_i}$  and the time delays  $\tau_i(t_k)$  for each path, are constant over time and known.

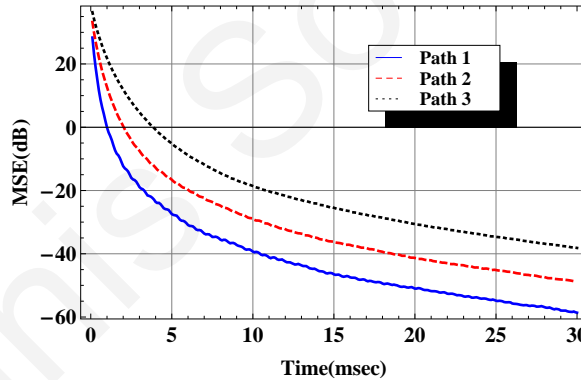
The performance of the derived estimators of Section 4.4.3 is evaluated using the Mean Square Error (4.77) of each estimator for each path over a period of time of 30msec, for  $N = 100$  realizations. The Signal to Noise ratio for the first receiving path is assumed to be 10dB. A smaller SNR is used compared to the previous experiments since here the real distributions of the parameters are known. This is the reference SNR for the system. The SNR of the second and third receiving paths is 9dB and 8dB respectively.

Starting with the phase estimation, Figure 4.3(a) displays the performance of the phase estimator (4.71) over time for each arriving path for a reference  $SNR = 10dB$ . The MSE decreases with time for all three paths, as time increases. This shows that the estimated phase converges in mean square sense to the real phase, as time increases, for all three paths.



**Figure 4.3:** MSE of the phase estimator

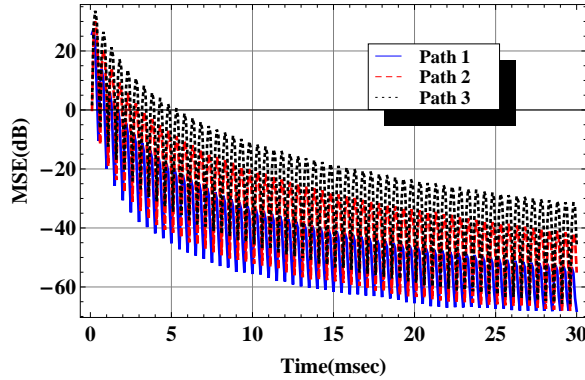
Similarly, the performance of the channel attenuation estimates is evaluated. Figure 4.4 displays the performance of the attenuation estimator (4.73) over time for each arriving path for a reference  $SNR = 10dB$ . The observations, are similar with the ones experienced for the phase MSE. Clearly, the MSE decreases as time increases.



**Figure 4.4:** MSE of the attenuation estimator, for a reference SNR=10 dB

Next, the conditional least-square estimate of the noiseless received signal for every path, defined by  $h_i(t_m, \theta_i, r_i) \triangleq E[h_i(t_m, \theta_i, r_i) | \mathcal{F}_m^y]$  is considered. Figure 4.5 displays the performance of the estimator (4.69) over time for each arriving path for a reference  $SNR = 10dB$ . The estimator converges to the real value, in mean square sense, as time increases for all three paths. This is verified by the fact that the MSE decrease to a very small value as time increases.

Finally, when comparing the results of the derived estimators it is obvious that the



**Figure 4.5:** MSE of the noiseless received signal estimator, for a reference SNR=10 dB

attenuation estimator and the noiseless received signal estimator perform better than the phase estimator. Since a  $10dB$  SNR can be considered as a low to medium system, the reference SNR is increased to  $20dB$  and the phase estimator is once again evaluated. Figure 4.3(b) displays the performance of the phase estimator (4.71) over time for each arriving path for a reference  $SNR = 20dB$ . Notice that the estimator produces better results, in MSE sense, than the previous evaluation. Clearly, SNR influences the performance of the estimator.

## 4.5 Summary

This chapter considers robust estimation for autoregressive channel models, and employs the KL distance to describe the uncertainty. The methodology presented addresses the maximization problem using variational methods, while for the minimization problem it invokes a change of probability measure technique. The change of probability measure technique introduced is being used in order to derive a recursive equation for the unnormalized  $\hat{a}$  posteriori distribution of nonlinear filtering problems. The theory developed is applied to two specific models, the linear Gaussian model and the non-coherent multipath model. The results derived include new robust least-square estimators for both estimation problems. Moreover, the robust minimax estimators for the non-coherent model are linked with classical estimators (e.g., by reducing the model uncertainty to zero).

# CHAPTER 5

## GENERALIZED MAP AND ML ESTIMATION

In this chapter two well-known estimation techniques, the Maximum  $\hat{A}$  Posteriori (MAP) technique and the Maximum Likelihood (ML) technique are being investigated. First, in Section 5.2 the MAP estimation technique is presented and a new generalized estimator is derived. Then, in Section 5.3 a different approach in deriving a ML estimator is presented. These generalized estimators deal with situations when the true distribution is unknown but belongs to a specific set described by relative entropy constraint. They are also obtained by modifying the uniform cost function of the classical MAP technique. The chapter ends with specific examples.

### 5.1 Introduction

In Bayesian statistics, the maximum  $\hat{a}$  posteriori estimate is defined as the mode of the posterior distribution, if posterior is unimodal. The MAP can be used to obtain a point estimate of an unobserved quantity on the basis of empirical data. It is closely related to Fisher's method of maximum likelihood, but employs an augmented optimization objective which incorporates a prior distribution over the quantity one wants to estimate. MAP estimation can therefore be seen as a regularization of ML estimation.

MAP estimates can be computed in several ways [53].

1. Analytically, when the mode(s) of the posterior distribution can be given in closed form. This is the case when conjugate priors are used.

2. Via numerical optimization such as the conjugate gradient method or Newton's method. This usually requires first or second derivatives, which have to be evaluated analytically or numerically.
3. Via a modification of an expectation-maximization algorithm. This does not require derivatives of the posterior density.
4. Via a Monte Carlo method using simulated annealing.

While MAP estimation is a limit of Bayes estimators (under the  $0 - 1$  cost function), it is not very representative of Bayesian methods in general. This is because MAP estimates are point estimates, whereas Bayesian methods are characterized by the use of distributions to summarize data and draw inferences: thus, Bayesian methods tend to report the posterior mean or median instead, together with credible intervals [1].

The classical maximum *a posteriori* estimation theory is developed by using a uniform cost function of  $0 - 1$ . In this chapter a new technique is developed, using the principles of the MAP estimation theory, with the only difference lying in the way that the cost function, is defined. In this chapter instead of a uniform cost function an exponential one is being used. The goal is to derive a generalized MAP estimator which includes as a special case the classical estimator, when the cost function is uniform.

For many observation models arising in practise it is not possible to apply the above technique, either because of intractability of the required analysis or because of the lack of a useful complete statistic. For such models, an alternative non Bayesian method is the maximum likelihood estimation technique. The idea behind maximum likelihood parameter estimation is to determine the parameters that maximize the probability (likelihood) of the sample data. From the statistical point of view, maximum likelihood estimation is regarded to be robust while it yields estimators with good statistical and convergence properties. In other words, ML estimation methods are versatile and apply to most models and to different types of data. In addition, they provide efficient methods for quantifying uncertainty through confidence bounds. Although the methodology for maximum likelihood estimation is simple, the implementation is mathematically intense.

For "large" samples, several results from central limit theorem are applicable yielding ML estimators which have the following properties [53]:

1. ML estimators are asymptotically normally distributed.
2. ML estimators are asymptotically “minimum variance.”
3. ML estimators are asymptotically unbiased (ML estimators are often biased, but the bias  $\rightarrow 0$  as  $n \rightarrow \infty$ ).

Maximum likelihood estimation represents the backbone of statistical estimation. It is based on deep theory, originally developed by R. A. Fisher. While beginning classes often focus on least squares estimation (“regression”); likelihood theory is the omnibus approach across the sciences, engineering and medicine. This chapter deal with the methodology used to derive a ML estimator and applies a new approach. The goal is to derive a generalized ML estimator using, like the MAP case, an exponential function.

## 5.2 Generalized Maximum $\hat{A}$ Posteriori Estimation

In this section the Maximum  $\hat{A}$  Posteriori method is being investigated. First the problem is defined by introducing the appropriate spaces and the general concept. Then the generalized estimator is derived by using the theory behind the classical estimator but by applying an exponential cost function.

### 5.2.1 Abstract Formulation

Suppose a measurable space  $(\Omega, \mathcal{F})$  is given on which the unobserved Random Variable (RV),  $X$  and the observed RV  $Y$  are defined, via  $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \Sigma_{\mathcal{X}})$ ,  $Y : (\Omega, \mathcal{F}) \rightarrow (\mathcal{Y}, \Sigma_{\mathcal{Y}})$ .

Thus  $\mathcal{X}$  is the space of the unobserved RV, and  $\mathcal{Y}$  is the space of the observed RV. The relation between the unobserved RV  $X$  and the observed RV  $Y$  is defined via a probabilistic mapping. In the Chapter 3 two different probabilistic mappings were used, the mapping  $\mu : \mathcal{X} \times \Sigma_{\mathcal{Y}} \rightarrow [0, 1]$  and the mapping  $\eta : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$ . These two mapping, as defined in Sections 3.2.1 and 3.2.3 are also going to be used in this chapter.

The objective is to estimate  $X$  by a function of the random variable  $Y$ . Let  $\hat{X} = \Phi(Y)$  denote the estimate of  $X$ . The previous chapters deal with the problem of Least-Square estimation. Here, a different estimation is investigated, which also belongs to the Bayesian group, the Maximum  $\hat{A}$  Posteriori estimation method.

The general theory is based on introducing a pay-off, and then minimizing the pay-off over  $x \in \mathcal{X}$ . As mentioned in previous paragraph, both Least-Square estimation and MAP estimation belong to the same set of estimation methods, the Bayesian Estimation. What distinguishes the different estimation methods is the determination of the appropriate pay-off function.

### 5.2.2 Derivation of Generalized Estimator

Let  $\ell : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$  be an  $\Sigma_{\mathcal{X}} \times \Sigma_{\mathcal{X}}$ -measurable function, and let  $s \in \mathfrak{R}$ . Then a new cost function is introduced defined by

$$C(x, \Phi(y)) = \begin{cases} e^{\frac{\ell(x)}{s}}, & \text{if } |x - \Phi(y)| > \Delta \\ 0, & \text{if } |x - \Phi(y)| \leq \Delta \end{cases} \quad (5.1)$$

where  $\Delta > 0$ . As shown in Fig. 5.1, the cost function  $C(x, \Phi(y))$  for  $|x - \Phi(y)| > \Delta$  is a convex increasing function of the error. The classical MAP estimation assumes a uniform cost function  $C(x, \Phi(y)) = 1$  for  $|x - \Phi(y)| > \Delta$  and 0 for  $|x - \Phi(y)| \leq \Delta$ . Here it is assumed that the cost function is not uniform but an exponential function of  $x$ . Note that the classical case is included in this cost function, since  $s \rightarrow \infty$ , implies  $C(x, \Phi(y)) = 1$  for  $|X - \Phi(y)| > \Delta$  (red graph in Fig. 5.1).

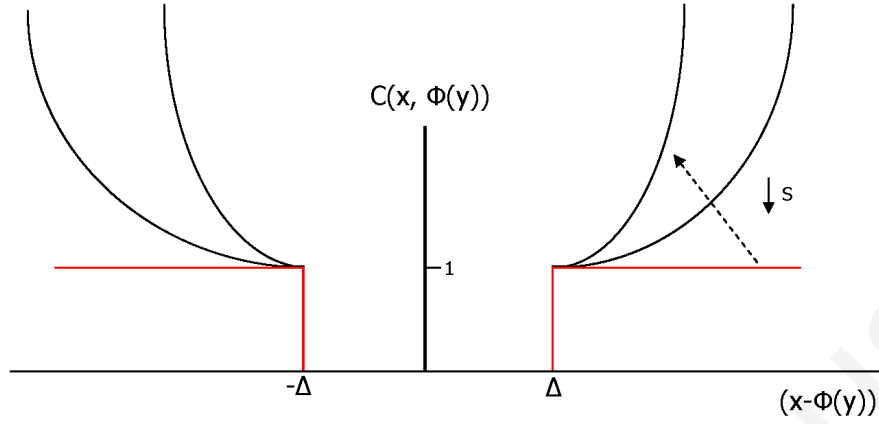
For the above cost function the average pay-off is given by

$$\begin{aligned} E[C(X, \Phi(Y))] &= \int_{\mathcal{X} \times \mathcal{Y}} e^{\frac{\ell(x)}{s}} I_{\{|x - \Phi(y)| > \Delta\}}(x) dP_{X,Y}(x, y) \\ &= \int_{\mathcal{Y}} dP_Y(y) \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} I_{\{|x - \Phi(y)| > \Delta\}}(x) \eta(y, dx) \end{aligned} \quad (5.2)$$

where  $I$  is an index function defined by

$$I_{\{|x - \Phi(y)| > \Delta\}}(w) = \begin{cases} 1, & \text{if } w \in \{x : |x - \Phi(y)| > \Delta\} \\ 0, & \text{if } w \notin \{x : |x - \Phi(y)| > \Delta\} \end{cases}$$




**Figure 5.1:** MAP Cost Function

Note that for a given  $Y = y$ , the space  $\mathcal{X}$  of the unobserved RV, is partitioned into

$$\mathcal{X} = \{x : |x - \Phi(y)| > \Delta\} \cup \{x : |x - \Phi(y)| \leq \Delta\} \quad (5.3)$$

and the index function  $I_{\mathcal{X}}$  of the space the unobserved RV is given by

$$I_{\mathcal{X}}(x) = I_{\{x:|x-\Phi(y)|>\Delta\}}(x) + I_{\{x:|x-\Phi(y)|\leq\Delta\}}(x) = 1. \quad (5.4)$$

Therefore, the index function  $I_{\{|x-\Phi(y)|>\Delta\}}(x)$  is given by

$$I_{\{x:|x-\Phi(y)|>\Delta\}}(x) = 1 - I_{\{x:|x-\Phi(y)|\leq\Delta\}}(x). \quad (5.5)$$

Using (5.5) with (5.2) the average pay-off is given by

$$\begin{aligned} E[C(X, \Phi(Y))] &= \int_{\mathcal{Y}} dP_Y(y) \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} \left[ 1 - I_{\{x:|x-\Phi(y)|\leq\Delta\}}(x) \right] \eta(y, dx) \\ &= \int_{\mathcal{Y}} dP_Y(y) \left[ \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} \eta(y, dx) - \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} I_{\{x:|x-\Phi(y)|\leq\Delta\}}(x) \eta(y, dx) \right]. \end{aligned} \quad (5.6)$$

Then, as mention before, the estimation is done by minimizing the average pay-off (5.6) over  $x \in \mathcal{X}$ . But, since  $dP_Y(y)$  is nonnegative and not a function of  $\Phi(y)$ , the estimation is done by minimizing the inner integral, which is given by

$$\begin{aligned} E[C(X, \Phi(Y))|Y = y] &= \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} I_{\{x:|x-\Phi(y)|>\Delta\}}(x) \eta(y, dx) \\ &= \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} \eta(y, dx) - \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} I_{\{x:|x-\Phi(y)|\leq\Delta\}}(x) \eta(y, dx). \end{aligned} \quad (5.7)$$

The first part of (5.7) given by  $\int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} \eta(y, dx)$  is a nonnegative number that can be estimated and which is also independent of  $x$ , so the minimization of the average pay-off (5.2) can be done by maximizing the second part of (5.7) given by

$$\begin{aligned} E \left[ e^{\frac{\ell(x)}{s}} I_{\{x: |x - \Phi(y)| > \Delta\}}(x) \right] &= \int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} I_{\{x: |x - \Phi(y)| \leq \Delta\}}(x) \eta(y, dx) \\ &= \int_{\Phi(y) - \Delta}^{\Phi(y) + \Delta} e^{\frac{\ell(x)}{s}} \eta(y, dx). \end{aligned} \quad (5.8)$$

If  $\eta(y, dx)$  is a smooth function of  $x$  and if  $\Delta = dx/2$  is sufficiently small, then

$$\int_{\Phi(y) - dx/2}^{\Phi(y) + dx/2} e^{\frac{\ell(x)}{s}} \eta(y, dx) \simeq e^{\frac{\ell(\Phi(y))}{s}} \eta(y, dx) \quad (5.9)$$

and the right-hand side is maximized by choosing  $\Phi(y)$  to be the value of  $x$  maximizing  $e^{\frac{\ell(x)}{s}} \eta(y, dx)$ , over  $x \in \mathcal{X}$ .

Therefore,

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}} e^{\frac{\ell(x)}{s}} \eta(y, dx). \quad (5.10)$$

However, if  $f_{\eta}(x|y) \triangleq \frac{\eta(y, dx)}{dx}$  exists then

$$\Phi_{MAP}(y) \triangleq \arg \max_{x \in \mathcal{X}} e^{\frac{\ell(x)}{s}} f_{\eta}(y|x). \quad (5.11)$$

In modeling a given statistical situation the family of conditional distributions (or stochastic kernels) of  $Y$  given  $X = x$  are needed, and for the Bayesian formulation the prior distribution for  $X$  is also needed. The conditional distribution of  $X$  given  $Y$  (or the stochastic kernel  $\eta(y, dx)$ ) can be obtained from the prior and the conditional of  $Y$  given  $X$  (or the stochastic kernel  $\mu(x, dy)$ ) by applying Bayes' formula.

$$\eta(y, dx) = \frac{\mu(x, dy) dP_X(x)}{\int_{\mathcal{X}} \mu(x, dy) dP_X(x)} = \frac{\mu(x, dy) dP_X(x)}{dP_Y(y)}. \quad (5.12)$$

The MAP estimator in (5.10) can be obtained using (5.12) but without the computation of  $dP_Y(y)$  since this term will not affect the maximization over  $x$ . That is,

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}} e^{\frac{\ell(x)}{s}} \mu(x, dy) dP_X(x). \quad (5.13)$$

Since the logarithm is an increasing function,  $\Phi_{MAP}(y)$  can be obtained by maximizing the following function

$$\log \left[ e^{\frac{\ell(x)}{s}} \frac{\mu(x, dy)}{dy} \frac{dP_X(x)}{dx} \right] = \frac{\ell(x)}{s} + \log \frac{\mu(x, dy)}{dy} + \log \frac{dP_X(x)}{dx}. \quad (5.14)$$

### 5.2.3 Connection with Minimax Approach

In the classical MAP estimation method, the MAP estimator  $\Phi_{MAP}(y)$  is obtained by choosing  $\Phi(y)$  to be the value of  $x$  maximizing the  $\acute{a}$  posteriori distribution  $\eta(y, dx)$ , over  $x \in \mathcal{X}$ . The connection of the new MAP estimate to robustness is described below.

It is assumed that this probabilistic kernel,  $\eta(y, dx)$ , represents the nominal system model or mapping ( $\acute{a}$  posteriori information) and that true kernel  $\nu(y, dx)$ , denoted by  $\nu : \mathcal{Y} \times \Sigma_{\mathcal{X}} \rightarrow [0, 1]$  belongs to an uncertainty set described by

$$\mathcal{B}(\eta) \triangleq \left\{ \nu \in \mathcal{P} : \int_{\mathcal{Y}} H(\nu|\eta)(y) dP_Y(y) \leq \int_{\mathcal{Y}} R(y) dP_Y(y) \stackrel{\nabla}{=} r_2 \right\},$$

similar with Section 3.2.3 of Chapter 3. The sample pay-off  $\ell(x, \Phi(y))$  used there is replaced by the pay-off  $\ell(x)$ , then the worst case measure of the true kernel is given by

$$\nu^*(y, dx) = \frac{e^{\frac{\ell(x)}{s}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x)}{s}} \eta(y, dx)} \quad (5.15)$$

and the MAP estimator  $\Phi_{MAP}(y)$  is obtained by choosing  $\Phi(y)$  to be the value of  $x$  maximizing  $\nu^*(y, dx)$ , over  $x \in \mathcal{X}$ . This is the same as maximizing only the numerator.

That is,

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}} e^{\frac{\ell(x)}{s}} \eta(y, dx), \quad (5.16)$$

which is the same result as obtained in Section 5.2.2 above.

## 5.3 Generalized Maximum Likelihood Estimation

In this section, the Maximum Likelihood estimation method is investigated. First the appropriate spaces are introduced and the problem is set up. Then the generalized estimator is derived by using the theory behind the classical estimator in combination with the theory of the generalized MAP estimator described in the previous section.

### 5.3.1 Abstract Formulation

For many observation models arising in practice, it is not possible to apply the Bayesian Estimation methods like Least-Square Estimation and Maximum  $\hat{A}$  Posteriori Estimation, mainly due to the lack of a useful complete sufficient statistic. For such models a very commonly used method of designing estimators is the Maximum Likelihood method.

Suppose that a measurable space  $(\Omega, \mathcal{F})$  is given on which the unobserved RV,  $X$  and the observed RV  $Y$  are defined, via  $X : (\Omega, \mathcal{F}) \rightarrow (\mathcal{X}, \Sigma_{\mathcal{X}})$ ,  $Y : (\Omega, \mathcal{F}) \rightarrow (\mathcal{Y}, \Sigma_{\mathcal{Y}})$ .

Thus  $\mathcal{X}$  is the space of the unobserved RV, and  $\mathcal{Y}$  is the space of the observed RV. The relation between the unobserved RV  $X$  and the observed RV  $Y$  is defined via the same probabilistic mapping as the previous section.

The classical Maximum Likelihood estimation technique can be derived from the Maximum  $\hat{A}$  Posteriori method. In the absence of any prior information about the RV  $X$  which has to be estimated, it is assumed that is uniformly distributed in its range since this represents more or less a worst-case prior. The same assumption is used in the generalized ML estimation technique described in this section.

### 5.3.2 Derivation of Generalized Estimator

The generalized MAP estimator derived in Section 5.2 is considered here, which is given by (5.13)

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}} e^{\frac{\ell(x)}{s}} \mu(x, dy) dP_X(x).$$

In the absence of any prior information about the parameter, it is assumed that it is uniformly distributed in its range since this represents more or less a worst-case prior. Applying the assumption above, the MAP estimate for a given  $y \in \mathcal{Y}$  is any value of  $x$  that maximizes  $e^{\frac{\ell(x)}{s}} \mu(x, dy)$  over  $x \in \mathcal{X}$ . Since  $\mu(x, dy)$  as a function of  $x$  is sometimes called the likelihood function, the classical estimate is called maximum likelihood estimate. Here the same name is used for the derived estimate, which is denoted by  $\Phi(y)_{ML}$  and is given by

$$\Phi_{ML}(y) = \arg \max_{x \in \mathcal{X}} e^{\frac{\ell(x)}{s}} \frac{\mu(x, dy)}{dy}. \quad (5.17)$$

Maximizing  $e^{\frac{\ell(x)}{s}} \mu(x, dy)$  is equivalent to maximizing  $\log \left[ e^{\frac{\ell(x)}{s}} \frac{\mu(x, dy)}{dy} \right] = \frac{\ell(x)}{s} + \log \frac{\mu(x, dy)}{dy}$ , and assuming sufficient smoothness of this function, a necessary condition for the maximum likelihood estimate is

$$\frac{\partial}{\partial x} \left[ \frac{\ell(x)}{s} + \log \frac{\mu(x, dy)}{dy} \right] \Big|_{x=\Phi_{ML}(y)} = 0. \quad (5.18)$$

Note that when  $s \rightarrow \infty$  then the generalized estimator converges to the classical ML estimator.

## 5.4 Examples

Next, some examples are presented to illustrate applicability of the theoretical results.

Given  $X$  is a RV and  $Y^m = \{Y_1, Y_2, \dots, Y_m\}$  is a sequence of RVs defined on  $(\Omega, \mathcal{F}, P)$ , which are related via the model

$$Y_i = HX + W_i, \quad i = 0, \dots, m \quad (5.19)$$

Hence  $Y^m : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}^{(m+1)d}, \mathcal{B}(\mathfrak{R}^{(m+1)d}))$  the observed sequence of RVs,  $X : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}^n, \mathcal{B}(\mathfrak{R}^n))$  the unobserved RV, and  $W^m : (\Omega, \mathcal{F}) \rightarrow (\mathfrak{R}^{(m+1)d}, \mathcal{B}(\mathfrak{R}^{(m+1)d}))$  the noise sequence of RVs. It is assumed  $W_i$  and  $X$  are independent Gaussian Random Variables,  $N(0, \Sigma_W)$ ,  $\Sigma_W > 0$ ,  $N(0, \Sigma_X)$ ,  $\Sigma_X > 0$ . Assuming (5.19) denotes the system model, then

$$\mu(x, dy^m) = \prod_{i=0}^m \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}} e^{-(y_i - Hx)^T \frac{\Sigma_W^{-1}}{2} (y_i - Hx)} dy_i, \quad (5.20)$$

$$P_X(dx) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_X|^{\frac{1}{2}}} e^{-x^T \frac{\Sigma_X^{-1}}{2} x} dx. \quad (5.21)$$

Also the following function is defined

$$\ell(x) = x^T U x, \quad U = U^T > 0. \quad (5.22)$$

Complete derivations of the results presented below can be found in Appendix E.

### 5.4.1 Generalized MAP Estimator

In this first example the generalized MAP estimator for the model (5.19) is computed together with the mean square error for this MAP estimator.

First the MAP estimator  $\Phi_{MAP}(y)$  is given by

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}} \left\{ \frac{\ell(x)}{s} + \log \frac{\mu(x, dy)}{dy} + \log \frac{dP_X(x)}{dx} \right\}. \quad (5.23)$$

Using (5.20), (5.21) and (5.22) the MAP estimator is given by

$$\Phi_{MAP}(y) = \arg \max_{x \in \mathcal{X}_{ad}} \bar{L}(x, y^m) \quad (5.24)$$

where

$$\begin{aligned} \bar{L}(x, y^m) &= \frac{x^T U x}{s} + \log \left( \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}} \right)^{m+1} - \sum_{i=0}^m (y_i - Hx)^T \frac{\Sigma_W^{-1}}{2} (y_i - Hx) \\ &+ \log \left( \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_X|^{\frac{1}{2}}} \right) - x^T \frac{\Sigma_X^{-1}}{2} x. \end{aligned} \quad (5.25)$$

The MAP estimator  $\Phi_{MAP}(y^m)$  of  $X$  is computed by differentiating the above pay-off function and setting the derivative to zero as shown below.

$$\frac{\partial}{\partial x} \bar{L}(x, y^m) \Big|_{x=\Phi_{MAP}} = 0. \quad (5.26)$$

So, by differentiating (5.25) the final generalized MAP estimator is derived

$$\Phi_{MAP}(y^m) = \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{s} \right)^{-1} H^T \Sigma_W^{-1} \sum_{i=0}^m y_i. \quad (5.27)$$

Next the mean square error for  $\Phi_{MAP}(y^m)$  is computed. The MSE is given by

$$\begin{aligned} &E \left[ (X - \Phi_{MAP}(Y^m))^T (X - \Phi_{MAP}(Y^m)) \right] = \\ &tr \left( E \left[ (X - \Phi_{MAP}(Y^m))(X - \Phi_{MAP}(Y^m))^T \right] \right). \end{aligned} \quad (5.28)$$

The mean square error can also be written as

$$\begin{aligned} &E \left[ (X - \Phi_{MAP}(Y^m))^T (X - \Phi_{MAP}(Y^m)) \right] = \\ &tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} (x - \Phi_{MAP}(y^m))(x - \Phi_{MAP}(y^m))^T \mu(x, dy) dP_X(x). \end{aligned} \quad (5.29)$$

After some manipulation the above function can be expressed as

$$\begin{aligned}
& E \left[ (X - \Phi_{MAP}(Y^m))^T (X - \Phi_{MAP}(Y^m)) \right] = \\
& tr \left( \Sigma_X - 2(m+1)H^T \Sigma_{\Phi_{MAP}}^T \Sigma_X + (m+1)^2 H^T \Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} H \Sigma_X \right. \\
& \left. + (m+1) \Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \Sigma_W \right) \tag{5.30}
\end{aligned}$$

where  $\Sigma_{\Phi_{MAP}} = \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{s} \right)^{-1} H^T \Sigma_W^{-1}$ .

### 5.4.2 Generalized ML Estimator

In the second example the generalized ML estimator for the model (5.19) is computed together with the mean square error for this ML estimator. The *Cramér – Rao lower bound* is also investigated. For this example it is assumed that  $X = x$  is a parameter.

The ML estimator  $\Phi_{ML}(y)$  is given by

$$\Phi_{ML}(y) = \arg \max_{x \in \mathcal{X}_{ad}} \left\{ \frac{\ell(x)}{s} + \log \frac{\mu(x, dy)}{dy} \right\}. \tag{5.31}$$

Using (5.20) and (5.22) the ML estimator is given by

$$\Phi_{ML}(y) = \arg \max_{x \in \mathcal{X}} \bar{L}_2(x, y^m) \tag{5.32}$$

where

$$\bar{L}_2(x, y^m) = \frac{x^T U x}{s} + \log \left( \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}} \right)^{(m+1)} - \sum_{i=0}^m (y_i - Hx)^T \frac{\Sigma_W^{-1}}{2} (y_i - Hx). \tag{5.33}$$

The generalized ML estimator  $\Phi_{ML}(y)$  of  $X$  is computed by differentiating the above pay-off function and setting the derivative to zero as shown below.

$$\frac{\partial}{\partial x} \bar{L}_2(x, y^m) \Big|_{x=\Phi_{ML}} = 0. \tag{5.34}$$

Therefore, by differentiating (5.33) the final generalized ML estimator is derived

$$\Phi_{ML}(y^m) = \left( (m+1)H^T \Sigma_W^{-1} H - \frac{2U}{s} \right)^{-1} H^T \Sigma_W^{-1} \sum_{i=0}^m y_i. \tag{5.35}$$

Next, the mean square error for  $\Phi_{ML}(y^m)$  is computed. The MSE is given by

$$\begin{aligned} & E\left[(x - \Phi_{ML}(Y^m))^T(x - \Phi_{ML}(Y^m))\right] \\ &= \text{tr}\left(E\left[(x - \Phi_{ML}(Y^m))(x - \Phi_{ML}(Y^m))^T\right]\right). \end{aligned} \quad (5.36)$$

Given the assumption that in ML estimation  $X$  is deterministic, after some manipulation the above function is expressed as

$$\begin{aligned} & E\left[(x - \Phi_{ML}(Y^m))^T(x - \Phi_{ML}(Y^m))\right] = \\ & \text{tr}\left(xx^T - 2(m+1)H^T\Sigma_{\Phi_{ML}}^T xx^T + (m+1)^2H^T\Sigma_{\Phi_{ML}}^T\Sigma_{\Phi_{ML}}Hxx^T\right. \\ & \left.+ (m+1)\Sigma_{\Phi_{ML}}^T\Sigma_{\Phi_{ML}}\Sigma_W\right) \end{aligned} \quad (5.37)$$

where  $\Sigma_{\Phi_{ML}} = \left((m+1)H^T\Sigma_W^{-1}H - \frac{2U}{s}\right)^{-1}H^T\Sigma_W^{-1}$ .

Finally, it is going to be investigated if this new generalized ML estimator satisfies the *Cramér – Rao lower bound*.

First, it is determined whether the estimator is biased, so the expectation of the estimator  $\Phi_{ML}(y^m)$  is computed.

$$E[\Phi_{ML}(Y^m)] = \Sigma_{\Phi_{ML}} \sum_{i=0}^m E[Y_i] = \Sigma_{\Phi_{ML}} \sum_{i=0}^m Hx = \Sigma_{\Phi_{ML}}(m+1)Hx \neq x. \quad (5.38)$$

This shows that the estimator is biased.

The *Cramér – Rao Inequality* for biased estimator is given by

$$\text{Var}[\Phi_{ML}(Y^m)] \geq V_X I_X^{-1} V_X^T \quad (5.39)$$

where

$$[V_X]_{ij} \triangleq \frac{\partial}{\partial x_j} E_X[\Phi_{ML}(y^m)_i], 1 \leq i, j \leq n \quad (5.40)$$

$$[I_X]_{ij} \triangleq -E_X\left[\frac{\partial^2}{\partial x_i \partial x_j} \log \frac{\mu(x, dy^m)}{dy^m}\right]. \quad (5.41)$$

Next the LHS of (5.39) is computed

$$\begin{aligned} \text{Var}[\Phi_{ML}(Y^m)] &= E\left[\left(\Phi_{ML}(Y^m) - E[\Phi_{ML}(Y^m)]\right)\left(\Phi_{ML}(Y^m) - E[\Phi_{ML}(Y^m)]\right)^T\right] \\ &= (m+1)\Sigma_{\Phi_{ML}}\Sigma_W\Sigma_{\Phi_{ML}}^T. \end{aligned} \quad (5.42)$$



Then the RHS of (5.39) is computed. Given (5.40),

$$V_X = (m + 1)\Sigma_{\Phi_{ML}}H \quad (5.43)$$

and given (5.41),

$$I_X^{-1} = \left( (m + 1)H^T\Sigma_w^{-1}H \right)^{-1} \quad (5.44)$$

provided that  $H$  is a square and invertible matrix.

Then

$$V_X I_X^{-1} V_X^T = (m + 1)\Sigma_{\Phi_{ML}}H \left( (m + 1)H^T\Sigma_w^{-1}H \right)^{-1} (m + 1)H^T\Sigma_{\Phi_{ML}}^T. \quad (5.45)$$

Therefore, assuming that  $H$  has size  $d \times d$  and is invertible, then

$$\begin{aligned} V_X I_X^{-1} V_X^T &= (m + 1)\Sigma_{\Phi_{ML}}H \left( (m + 1)H^T\Sigma_w^{-1}H \right)^{-1} (m + 1)H^T\Sigma_{\Phi_{ML}}^T \\ &= (m + 1)\Sigma_{\Phi_{ML}}HH^{-1}\Sigma_w(H^T)^{-1}H^T\Sigma_{\Phi_{ML}}^T = (m + 1)\Sigma_{\Phi_{ML}}\Sigma_W\Sigma_{\Phi_{ML}}^T. \end{aligned} \quad (5.46)$$

Given the above, the *Cramér – Rao lower bound* is satisfied with equality

$$\text{Var}[\Phi_{ML}(Y^m)] = V_X I_X^{-1} V_X^T. \quad (5.47)$$

## 5.5 Summary

This chapter considers a generalized approach to the classical MAP and ML estimation. A new cost function is used for the case of MAP estimation and a generalized estimator is derived. Furthermore, a new approach is used, closely related to the MAP theory developed, in order to derive a generalized estimator for the ML estimation problem. The theory developed is applied to examples in order to illustrate the applicability of the results.



# CHAPTER 6

## CONCLUSION

In this chapter a summary of the main results of the thesis is presented alongside with some suggestions for future work plans.

### 6.1 Synopsis

This thesis deals with the subject of signal estimation from noisy measurements under uncertainty. The work presented is focused mainly on the problem of robust least-square estimation for uncertain systems. In addition, it also investigates the methodology behind two well-known and used estimation techniques, the Maximum  $\hat{A}$  Posteriori estimation technique and the Maximum Likelihood estimation technique.

Robust estimation is a common problem appearing frequently in statistics and signal processing and its been under study for many years. Usually the performance of a system depends on the  $\hat{a}$  priori knowledge of its input characteristics and even small deviations from the assumed conditions can affect it drastically. One of the most widely used ways on dealing with modeling uncertainties is the minimax approach. This approach has been applied to many detection and estimation problems and has been used for designing robust schemes by optimizing the worst case performance. This thesis uses probability distributions, or general measures, defined on measurable spaces to model the the uncertainty description and the nominal description of the system. Specifically, two type of uncertainty models are being considered: 1) Uncertainty Models on Conditional Distributions or otherwise known Stochastic

Kernels; 2) Uncertainty Models on Joint Distributions.

Stochastic kernel uncertainty models are being implemented in the design of communication systems, where the input message has a known distribution, while the channel is unknown but belongs to a certain class of channels. Here two types of uncertain stochastic models are being considered: i) when the conditional probability distribution of the measurement  $Y$  given the signal to be estimated  $X$ , or channel kernel, is unknown; ii) when the *a posteriori* distribution of  $X$  given  $Y$  is unknown. Joint distribution uncertainty models, on the other hand, are usually employed when both the unobserved and observed random variables are uncertain. The uncertainty description of these models is characterized by the class of uncertain measures which satisfy a relative entropy constraint with respect to a nominal measure. The problem of robust estimation is formulated by minimizing over the set of estimators, the maximum of a linear functional of the uncertain measure over the Kullback-Leibler (KL) distance constraint set.

In this thesis a general framework is presented where the basic ideas are explained and the fundamental results are derived. Abstract Polish spaces are being used to describe the least-square estimation problem, while, as was mentioned above, the uncertain model considered is described by stochastic kernels and joint distributions. Once the appropriate space of measures are introduced the maximizing kernels and joint measure are computed explicitly using Lagrangian functionals and variational methods. Furthermore, certain properties of the optimal solution are presented, including performance bounds, which are important for the numerical computation of the optimal solution. The theoretical results obtained are applied to various examples. The theory developed is also applied to Multiple-Input Multiple-Output (MIMO) communication systems, and to finite-dimensional autoregressive channel models, in order to derive robust estimators for a class of uncertain models. The methodology is presented, where the maximization is addressed using variational methods, while the minimization is addressed using a change of probability measure technique. The change of probability measure technique introduced is being used in order to derive recursive equations for the conditional distribution of nonlinear filtering problem. The theory is then applied first to a linear Gaussian model and then to an attenuated sinusoid in a multipath environment, which is subject to an additive Gaussian noise, i.e., to a non-coherent multipath model. For the linear Gaussian model a robust estimator is derived, which resembles the Kalman filter. For the non-coherent multipath model robust phase and enveloped estimators are

derived, and a connection with the classical least-square estimation is presented.

Finally, in addition to the robust least-square estimation, this thesis takes a closer look into the derivation of estimators for the Maximum  $\hat{A}$  Posteriori estimation method and the Maximum Likelihood estimation method. The MAP estimation technique is based on the minimization of a specific uniform cost function. In this thesis, the classical cost function is replaced by an exponential cost function, and a generalized MAP estimator is derived. A connection with the minimax approach used in the least-square problem is also presented. This generalized estimator includes as a special case the classical MAP estimator. A similar approach is also implemented for the case of the ML estimation technique. The methodology used for ML estimation involves the maximization of a likelihood function. Starting with the MAP estimation method and assuming that the parameter to be estimated is uniformly distributed in its range, the new likelihood function is derived. Through the above approach a generalized likelihood function is introduced, which also includes an exponential function. Examples are presented that illustrate the application of the derived results to various problems.

## 6.2 Directions for Future Research

The robust least-square estimation method together with the new generalized MAP and ML estimators presented in this thesis raise a few issues related to

1. Robust signal detection;
2. Power Spectral Density uncertainty models;
3. Total variation distance constraint between measures;
4. Application of the models to real-time applications, which involve also MIMO applications;
5. Information theory applications;
6. Relation of Mutual Information and MMSE.

This thesis deals with the subject of signal estimation. Besides signal estimation, one of the most pervasive of functions that signal processing schemes are required to carry out is that of detecting a signal of a generally known type in noisy observations. Most signal detection problems can be cast in the framework of  $M$ -ary hypothesis testing, in which there is an observation, usually a vector or function, on the basis of which a decision has to be made among  $M$  possible statistical situations describing the observation [1]. Obvious examples of applications in which signal detection is required are provided by radar (detection of echo pulses) and sonar (detection of a random signal present in an array of hydrophones). Numerous other applications may be listed; for example, detection of specified two-level pulse-code sequences in communication systems, and detection of abnormal patterns in medical imaging. The subject of signal detection and estimation deals with the processing of information-bearing signals in order to make inferences about the information that they contain.

The most fundamental problem of signal detection is the determination of the likelihood ratio for detecting signals against a noise background. This is a problem with a rich history in both electrical engineering and mathematics, and its solution has involved, over the years, a wide variety of mathematical tools and flashes of intuition. For one thing, mathematical models are often significant simplifications/idealizations of complex physical problems. Secondly, even if the model is reasonably good, the knowledge of the parameters in it, e.g., covariance functions, time constants, etc., may not be enough to justify a direct numerical evaluation of formulas derived from the model. The major engineering goal is to obtain structural insights into the mathematical solutions of classes of special problems, with the hope that these insights can then be used to intelligently modify and adapt the mathematical solution to the particular physical problem at hand [4]. The work presented in this thesis can be used in detection problems in order to obtain robust detection scheme. The uncertainty models can be implemented to describe the systems. Also, as was shown in Section 4.2.3 a change of probability measure can be used to obtain the conditional distribution. In decision applications, there will be one distribution  $\alpha_m^s(x)$  for each hypothesis, e.g., in binary hypothesis the decision rule takes the form

$$\frac{\inf_{\tilde{x}_m \in \mathcal{X}_{ad}} \int_{\mathbb{R}^n} d\alpha_m^{s,1}(x)}{\inf_{\tilde{x}_m \in \mathcal{X}_{ad}} \int_{\mathbb{R}^n} d\alpha_m^{s,2}(x)} \underset{H_2}{\overset{H_1}{\gtrless}} \gamma$$

where  $\gamma$  is the threshold.

Moreover, another area for future research, which is related to signal detection problems, is to devise decision rules for a class of probability distribution.

In statistical signal processing and physics, the spectral density, power spectral density is a positive real function of a frequency variable associated with a stationary stochastic process, or a deterministic function of time, which has dimensions of power per Hz. It is often called simply the spectrum of the signal. Intuitively, the spectral density captures the frequency content of a stochastic process and helps identify periodicities.

The concept and use of the power spectrum of a signal is fundamental in electronic engineering, especially in electronic communication systems (radio & microwave communications, radars, and related systems). As was mentioned in Chapter 1, the classes of allowable characteristics one deals with in robust signal processing are generally nonparametric function classes, such as the class of all power spectral density functions with specified total power (area under the function) and which lie between specified upper and lower bounding functions. Uncertainty models using the concept of PSD for robust estimation problems can be found in [15] and [30]. It will be interesting to combine the minimax approach presented in this thesis with the concept of power spectrum density and also use the concept of KL distance to define the constraint sets. The KL distance has been introduced between spectral density functions of stationary stochastic processes in [54].

In this thesis the uncertainty description of the system, and the nominal description of the system are modeled by probability distributions, or general measures, defined on measurable spaces. The uncertainty description of these models is characterized by the class of uncertain measures which satisfy a KL distance constraint with respect to a nominal measure. Over the last few years, the KL distance uncertainty model has received particular attention due to various properties (convexity, compact level sets), its simplicity and its connection to risk sensitive pay-off, minimax games, and large deviations. Unfortunately, KL distance uncertainty modeling has two disadvantages: 1) it does not define a true metric on the space of measures; 2) relative entropy between two measures is not defined if the measures are not absolutely continuous. The latter rules out the possibility of measures  $\nu \in \mathcal{M}_1(\Sigma)$  and  $\mu \in \mathcal{M}_1(\Sigma)$  to be defined on different spaces<sup>1</sup>. It is one of the main disadvantages

---

<sup>1</sup>This corresponds to the case in which the nominal system is a simplified version of the true system and is defined on a lower dimension space.

of employing relative entropy in the context of uncertainty modeling for stochastic controlled diffusions (or stochastic differential equations). Motivated by the above issues, the KL distance constrain can be replaced by an uncertainty model based on the total variation distance defined on the space measures. This uncertainty set is described by a ball with respect to the total variation norm, centered at the nominal measure having positive radius.

Given a known or nominal probability measure  $\mu \in \mathcal{M}_1(\Sigma)$  the uncertainty set based on total variation distance is defined by

$$B_R(\mu) \triangleq \left\{ \nu \in \mathcal{M}_1(\Sigma) : \|\nu - \mu\| \leq R \right\}$$

where  $R \in [0, \infty)$ . The total variation distance<sup>2</sup> on  $\mathcal{M}_1(\Sigma) \times \mathcal{M}_1(\Sigma)$  is defined by

$$\|\alpha - \beta\| \triangleq \sup_{P \in \mathcal{P}(\Sigma)} \sum_{F_i \in P} |\alpha(F_i) - \beta(F_i)|, \quad \alpha, \beta \in \mathcal{M}_1(\Sigma)$$

where  $\mathcal{P}(\Sigma)$  denotes the collection of all finite partitions of  $\Sigma$ . The above distance satisfies the properties of a metric, and does not require absolute continuity of measures when defining the uncertainty ball, i.e., singular measures are admissible and the measures need not be defined on the same space. It can very well be the case that  $\mu \in \mathcal{M}_1(\Sigma)$  and  $\nu \in \mathcal{M}_1(\tilde{\Sigma})$  where  $\Sigma \subset \tilde{\Sigma}$ . Since  $\mathcal{M}_1(\Sigma)$  are probability measures then it follows that the radius of uncertainty belongs to the restricted set  $R \in [0, 2]$ .

Clearly, the total variation distance uncertainty set is larger than the KL distance uncertainty set. This can be concluded from Pinsker's inequality [55] as follows.

$$\|\nu - \mu\|^2 \leq 2H(\nu|\mu), \quad \nu, \mu \in \mathcal{M}_1(\Sigma), \quad \text{if } \nu \ll \mu, \quad \log \frac{d\nu}{d\mu} \in L_1(\nu).$$

Hence, even for those measures which satisfy  $\nu \ll \mu$  and  $\log \frac{d\nu}{d\mu} \in L_1(\nu)$  the uncertainty set described by relative entropy is a subset of the much larger total variation distance uncertainty set, that is,  $A_{\frac{R^2}{2}}(\mu) \subset B_R(\mu)$ . In the parlance of stochastic differential equations the total variation distance covers the case when both drift and diffusion coefficients of the stochastic differential equations are uncertain.

This thesis presents mostly theoretical. A natural extension will be to apply the results to real-time applications and evaluate their performance. Special emphasis could be given to MIMO communication systems as they are vastly used in real time.

<sup>2</sup>The definition of total variation distance applies to signed measures as well.



Finally, it is noted that aside from minimax estimation, the tools developed in this thesis can also be used in Information theory. For example, in computing maxmin capacity for a class of channels, and minimax rate distortion for a class of sources.

The mutual information, which is at the core of information theory, is an indicator of how much coded information can be pumped through a channel reliably given a certain input signaling. As was already presented in Section 1.1.4, [31] presents a new formula that connects the input-output mutual information and the minimum mean-square error achievable by optimal estimation of the input given the output. That is, given that the input-output mutual information and the MMSE are monotone functions of the signal-to-noise ratio (SNR), denoted by  $I(\text{snr})$  and  $\text{mmse}(\text{snr})$ , respectively, the mutual information in nats and the MMSE satisfy the following relationship regardless of the input statistics:

$$\frac{d}{d\text{snr}} I(\text{snr}) = \frac{1}{2} \text{mmse}(\text{snr}). \quad (6.1)$$

This relationship holds for both scalar and vector signals. Also based on this relationship various other properties of MMSE are investigated in [32] and [33], like the MMSE dimension. All the above results assume that there is no uncertainty in the model being used. The work presented in this thesis, can be used to developed a new relationship between mutual information and MMSE for a class of models which are subject to uncertainty.



# APPENDIX A

## BASIC MATRIX IDENTITIES

In this thesis the following Matrix identities are being used [56].

1. Basic Identity. Given matrices  $A$  and  $B$  (both size  $n \times n$  and invertible), then

$$(AB)^{-1} = B^{-1}A^{-1}. \quad (\text{A.1})$$

2. Woodbury Identity. Given invertible matrix  $A$  (size  $n \times n$ ), matrix  $C$  (size  $n \times d$ ), and invertible matrix  $B$  (size  $d \times d$ ), then

$$(A + CBC^T)^{-1} = A^{-1} - A^{-1}C(B^{-1} + C^T A^{-1}C)^{-1}C^T A^{-1}. \quad (\text{A.2})$$

3. PosDef Identity. If matrices  $P$  and  $R$  are assumed to be positive definite and invertible, then

$$(P^{-1} + B^T R^{-1}B)^{-1}B^T R^{-1} = PB^T(BPB^T + R)^{-1}. \quad (\text{A.3})$$

Moreover, the following identities for Determinants are being used in this thesis.

- I. Given square matrices  $A$  and  $B$ , then

$$|A^T| = |A|, \quad (\text{A.4})$$

$$|AB| = |A||B|, \quad (\text{A.5})$$

$$|A^n| = |A|^n. \quad (\text{A.6})$$

- II. Given a square and invertible matrix  $A$ , then

$$|A|^{-1} = \frac{1}{|A|}. \quad (\text{A.7})$$

III. Given matrix  $A$  with size  $n \times m$  and matrix  $B$  with size  $m \times n$ , then

$$|I_n + AB| = |I_m + BA| \quad (\text{A.8})$$

where matrix  $I_n$  has size  $n \times n$  and matrix  $I_m$  has size  $m \times m$ .

Yiannis Socratous

## APPENDIX B

### PROOF OF REMARK 2.2.20

According to Theorem 2.2.19,  $\{\bar{\alpha}_k(x)\} \triangleq \{\pi_k(x)\}$  satisfies the following recursion

$$\bar{\alpha}_k(x) = \frac{\Xi_{v_k}(D_k^{-1}(y_k - C_k x))}{|D_k| \Xi_{v_k}(y_k)} \int_{\mathfrak{R}_n} \frac{\Psi_{w_k}(B_k^{-1}(x - A_k z))}{|B_k|} \bar{\alpha}_{k-1}(z) dz \quad (\text{B.1})$$

where

$$\Xi_{v_k}(D_k^{-1}(y_k - C_k x)) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{(y_k - C_k x)^T (D_k D_k^T)^{-1} (y_k - C_k x)}{2}\right), \quad (\text{B.2})$$

$$\Xi_{v_k}(y_k) = \frac{1}{(2\pi)^{\frac{d}{2}}} \exp\left(-\frac{y_k^T y_k}{2}\right), \quad (\text{B.3})$$

$$\Psi_{w_k}(B_k^{-1}(x - A_k z)) = \frac{1}{(2\pi)^{\frac{n}{2}}} \exp\left(\frac{-(x - A_k z)^T (B_k B_k^T)^{-1} (x - A_k z)}{2}\right). \quad (\text{B.4})$$

The solution of (B.1) is assumed to have the following form.

$$\bar{\alpha}_k(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |V_{k|k}|^{\frac{1}{2}}} \exp\left(- (x - \hat{x}_{k|k})^T \frac{(V_{k|k})^{-1}}{2} (x - \hat{x}_{k|k}) + \beta_{k|k}\right). \quad (\text{B.5})$$

Then,

$$\begin{aligned} \bar{\alpha}_{k-1}(z) &= \frac{1}{(2\pi)^{\frac{n}{2}} |V_{k-1|k-1}|^{\frac{1}{2}}} \\ &\times \exp\left(- (z - \hat{x}_{k-1|k-1})^T \frac{(V_{k-1|k-1})^{-1}}{2} (z - \hat{x}_{k-1|k-1}) + \beta_{k-1|k-1}\right). \end{aligned} \quad (\text{B.6})$$

The derivation starts with the integral part of (B.1). Then using (B.4) and (B.6) this can be expressed as

$$\int_{\mathfrak{R}_n} \frac{\Psi_{w_k}(B_k^{-1}(x - A_k z))}{|B_k|} \bar{\alpha}_{k-1}(z) dz = \int_{\mathfrak{R}_n} \frac{1}{(2\pi)^{\frac{n}{2}} |B_k| (2\pi)^{\frac{n}{2}} |V_{k-1|k-1}|^{\frac{1}{2}}}$$

$$\begin{aligned}
 & \times \exp \left( \frac{-(x - A_k z)^T (B_k B_k^T)^{-1} (x - A_k z)}{2} \right. \\
 & \left. - (z - \hat{x}_{k-1|k-1})^T \frac{(V_{k-1|k-1})^{-1}}{2} (z - \hat{x}_{k-1|k-1}) \right. \\
 & \left. + \beta_{k-1|k-1} \right). \tag{B.7}
 \end{aligned}$$

Next, the exponential term of B.7 is expanded as shown below.

$$\begin{aligned}
 & \exp \left( \frac{-(x - A_k z)^T (B_k B_k^T)^{-1} (x - A_k z)}{2} - (z - \hat{x}_{k-1|k-1})^T \frac{(V_{k-1|k-1})^{-1}}{2} (z - \hat{x}_{k-1|k-1}) \right. \\
 & \left. + \beta_{k-1|k-1} \right) \\
 & = -x^T \frac{(B_k B_k^T)^{-1}}{2} x - z^T A_k^T \frac{(B_k B_k^T)^{-1}}{2} A_k z + z^T A_k^T (B_k B_k^T)^{-1} x - z^T \frac{(V_{k-1|k-1})^{-1}}{2} z \\
 & \quad - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + z^T (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1} \\
 & = -\frac{1}{2} z^T \left( (V_{k-1|k-1})^{-1} + A_k^T (B_k B_k^T)^{-1} A_k \right) z \\
 & \quad + z^T \left( A_k^T (B_k B_k^T)^{-1} x + (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \right) \\
 & \quad - x^T \frac{(B_k B_k^T)^{-1}}{2} x - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1}.
 \end{aligned}$$

Now, the above equation can be written as a quadratic function of  $z$

$$\begin{aligned}
 = & \quad -(z - N_1)^T \frac{\Sigma_1^{-1}}{2} (z - N_1) + N_1^T \frac{\Sigma_1^{-1}}{2} N_1 - x^T \frac{(B_k B_k^T)^{-1}}{2} x \\
 & \quad - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1}
 \end{aligned}$$

where

$$\begin{aligned}
 \Sigma_1^{-1} &= (V_{k-1|k-1})^{-1} + A_k^T (B_k B_k^T)^{-1} A_k, \tag{B.8} \\
 \Sigma_1^{-1} N_1 &= A_k^T (B_k B_k^T)^{-1} x + (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \\
 \Rightarrow N_1 &= \left( (V_{k-1|k-1})^{-1} + A_k^T (B_k B_k^T)^{-1} A_k \right)^{-1} \\
 & \quad \times \left[ A_k^T (B_k B_k^T)^{-1} x + (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \right]. \tag{B.9}
 \end{aligned}$$

Using the fact that  $\Sigma_1^T = \Sigma_1$ ,  $(B_k B_k^T)^T = (B_k B_k^T)$  and  $(V_{k-1|k-1})^T = (V_{k-1|k-1})$ , then

$$\begin{aligned}
 N_1^T \frac{\Sigma_1^{-1}}{2} N_1 &= x^T (B_k B_k^T)^{-1} A_k \frac{\Sigma_1}{2} A_k^T (B_k B_k^T)^{-1} x \\
 & \quad + \hat{x}_{k-1|k-1}^T (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \\
 & \quad + x^T (B_k B_k^T)^{-1} A_k \Sigma_1 (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1}. \tag{B.10}
 \end{aligned}$$

Since

$$\int_{\mathbb{R}^n} \frac{1}{(2\pi)^{n/2} |\Sigma_1|^{1/2}} \exp\left(- (z - N_1)^T \frac{\Sigma_1^{-1}}{2} (z - N_1)\right) dz = 1$$

the integral term of (B.1) can be written as

$$\begin{aligned} & \int_{\mathbb{R}^n} \frac{\Psi_{w_k}(B_k^{-1}(x - A_k z))}{|B_k|} \bar{\alpha}_{k-1}(z) dz = \frac{|\Sigma_1|^{1/2}}{(2\pi)^{n/2} |B_k| |V_{k-1|k-1}|^{1/2}} \\ & \times \exp\left(N_1^T \frac{\Sigma_1^{-1}}{2} N_1 - x^T \frac{(B_k B_k^T)^{-1}}{2} x - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1}\right) \\ & \times \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{n/2} |\Sigma_1|^{1/2}} \exp\left(- (z - N_1)^T \frac{\Sigma_1^{-1}}{2} (z - N_1)\right) dz \\ & = \frac{|\Sigma_1|^{1/2}}{(2\pi)^{n/2} |B_k| |V_{k-1|k-1}|^{1/2}} \\ & \times \exp\left(N_1^T \frac{\Sigma_1^{-1}}{2} N_1 - x^T \frac{(B_k B_k^T)^{-1}}{2} x - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1}\right). \end{aligned} \quad (\text{B.11})$$

Next (B.11) is substituted into (B.1), using also (B.2) and (B.3), resulting in

$$\begin{aligned} \bar{\alpha}_k(x) &= \frac{|\Sigma_1|^{1/2}}{(2\pi)^{n/2} |D_k| |B_k| |V_{k-1|k-1}|^{1/2}} \exp\left(- (y_k - C_k x)^T \frac{(D_k D_k^T)^{-1}}{2} (y_k - C_k x)\right. \\ & \quad \left. + \frac{y_k^T y_k}{2} + N_1^T \frac{\Sigma_1^{-1}}{2} N_1 - x^T \frac{(B_k B_k^T)^{-1}}{2} x\right. \\ & \quad \left. - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1}\right) \end{aligned} \quad (\text{B.12})$$

The above expression can be simplified even more. First the exponential of (B.12) is simplified as shown below.

$$\begin{aligned} & \exp\left(- (y_k - C_k x)^T \frac{(D_k D_k^T)^{-1}}{2} (y_k - C_k x) + \frac{y_k^T y_k}{2} + N_1^T \frac{\Sigma_1^{-1}}{2} N_1 - x^T \frac{(B_k B_k^T)^{-1}}{2} x\right. \\ & \quad \left. - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1}\right) \\ &= -y_k^T \frac{(D_k D_k^T)^{-1}}{2} y_k - x^T C_k^T \frac{(D_k D_k^T)^{-1}}{2} C_k x + x^T C_k^T (D_k D_k^T)^{-1} y_k + \frac{y_k^T y_k}{2} \\ & \quad + x^T (B_k B_k^T)^{-1} A_k \frac{\Sigma_1}{2} A_k^T (B_k B_k^T)^{-1} x + \hat{x}_{k-1|k-1}^T (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \\ & \quad + x^T (B_k B_k^T)^{-1} A_k \Sigma_1 (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} - x^T \frac{(B_k B_k^T)^{-1}}{2} x \\ & \quad - \hat{x}_{k-1|k-1}^T \frac{(V_{k-1|k-1})^{-1}}{2} \hat{x}_{k-1|k-1} + \beta_{k-1|k-1} \end{aligned}$$

$$\begin{aligned}
 &= -\frac{1}{2}x^T \left[ C_k^T (D_k D_k^T)^{-1} C_k - (B_k B_k^T)^{-1} A_k \frac{\Sigma_1}{2} A_k^T (B_k B_k^T)^{-1} + (B_k B_k^T)^{-1} \right] x \\
 &+ x^T \left[ C_k^T (D_k D_k^T)^{-1} y_k + (B_k B_k^T)^{-1} A_k \Sigma_1 (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \right] \\
 &+ \hat{x}_{k-1|k-1}^T \left[ (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} - \frac{(V_{k-1|k-1})^{-1}}{2} \right] \hat{x}_{k-1|k-1} \\
 &+ y_k^T \left[ \frac{1}{2} I_d - \frac{(D_k D_k^T)^{-1}}{2} \right] y_k + \beta_{k-1|k-1}.
 \end{aligned}$$

Next, the above equation can be written as a quadratic function of  $x$

$$\begin{aligned}
 &= -(x - \hat{x}_{k|k})^T \frac{(V_{k|k})^{-1}}{2} (x - \hat{x}_{k|k}) \\
 &+ \hat{x}_{k|k}^T \frac{(V_{k|k})^{-1}}{2} \hat{x}_{k|k} + \hat{x}_{k-1|k-1}^T \left[ (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} - \frac{(V_{k-1|k-1})^{-1}}{2} \right] \hat{x}_{k-1|k-1} \\
 &+ y_k^T \left[ \frac{1}{2} I_d - \frac{(D_k D_k^T)^{-1}}{2} \right] y_k + \beta_{k-1|k-1} \tag{B.13}
 \end{aligned}$$

where

$$(V_{k|k})^{-1} = C_k^T (D_k D_k^T)^{-1} C_k - (B_k B_k^T)^{-1} A_k \Sigma_1 A_k^T (B_k B_k^T)^{-1} + (B_k B_k^T)^{-1} \tag{B.14}$$

$$\begin{aligned}
 (V_{k|k})^{-1} \hat{x}_{k|k} &= C_k^T (D_k D_k^T)^{-1} y_k + (B_k B_k^T)^{-1} A_k \Sigma_1 (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \\
 \Rightarrow \hat{x}_{k|k} &= (V_{k|k}) \left[ C_k^T (D_k D_k^T)^{-1} y_k + (B_k B_k^T)^{-1} A_k \Sigma_1 (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \right]. \tag{B.15}
 \end{aligned}$$

Using (B.8) with (B.14) then

$$\begin{aligned}
 (V_{k|k})^{-1} &= C_k^T (D_k D_k^T)^{-1} C_k - (B_k B_k^T)^{-1} A_k \left( (V_{k-1|k-1})^{-1} + A_k^T (B_k B_k^T)^{-1} A_k \right)^{-1} \\
 &\quad \times A_k^T (B_k B_k^T)^{-1} + (B_k B_k^T)^{-1}. \tag{B.16}
 \end{aligned}$$

Applying the Woodbury Identity (A.2), the above transforms to

$$(V_{k|k})^{-1} = C_k^T (D_k D_k^T)^{-1} C_k + \left( (B_k B_k^T) + A_k V_{k-1|k-1} A_k^T \right)^{-1} \tag{B.17}$$

$$(V_{k|k}) = \left( C_k^T (D_k D_k^T)^{-1} C_k + V_{k|k-1}^{-1} \right)^{-1} \tag{B.18}$$

where

$$\begin{aligned}
 V_{k|k-1} &= (B_k B_k^T) + A_k V_{k-1|k-1} A_k^T \tag{B.19} \\
 &= (B_k B_k^T) + A_k \left( C_{k-1}^T (D_{k-1}^T D_{k-1})^{-1} C_{k-1} + V_{k-1|k-2}^{-1} \right)^{-1} A_k^T
 \end{aligned}$$



$$\begin{aligned}
&= (B_k B_k^T) + A_k \left( V_{k-1|k-2} - V_{k-1|k-2} C_{k-1}^T (D_{k-1}^T D_{k-1} + C_{k-1} V_{k-1|k-2} C_{k-1}^T)^{-1} \right. \\
&\quad \left. C_{k-1} V_{k-1|k-2} \right) A_k^T \\
&= (B_k B_k^T) + A_k V_{k-1|k-2} A_k^T - A_k V_{k-1|k-2} C_{k-1}^T (D_{k-1}^T D_{k-1} \\
&\quad + C_{k-1} V_{k-1|k-2} C_{k-1}^T)^{-1} C_{k-1} V_{k-1|k-2} A_k^T.
\end{aligned} \tag{B.20}$$

Using (B.8) with (B.15), then

$$\begin{aligned}
\hat{x}_{k|k} &= (V_{k|k}) \left[ C_k^T (D_k D_k^T)^{-1} y_k \right. \\
&\quad \left. + (B_k B_k^T)^{-1} A_k \left( (V_{k-1|k-1})^{-1} + A_k^T (B_k B_k^T)^{-1} A_k \right)^{-1} (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \right].
\end{aligned} \tag{B.21}$$

Applying the PosDef Identity (A.3)

$$\begin{aligned}
\hat{x}_{k|k} &= (V_{k|k}) \left[ C_k^T (D_k D_k^T)^{-1} y_k \right. \\
&\quad \left. + \left( (B_k B_k^T) + A_k V_{k-1|k-1} A_k^T \right)^{-1} A_k V_{k-1|k-1} (V_{k-1|k-1})^{-1} \hat{x}_{k-1|k-1} \right] \\
&= (V_{k|k}) \left[ C_k^T (D_k D_k^T)^{-1} y_k + \left( (B_k B_k^T) + A_k V_{k-1|k-1} A_k^T \right)^{-1} A_k \hat{x}_{k-1|k-1} \right] \\
&= (V_{k|k}) \left[ C_k^T (D_k D_k^T)^{-1} y_k + (V_{k|k-1})^{-1} \hat{x}_{k|k-1} \right]
\end{aligned} \tag{B.22}$$

where

$$\hat{x}_{k|k-1} = A_k \hat{x}_{k-1|k-1}. \tag{B.23}$$

The above equation can be manipulated a bit more in order to get the following

$$\begin{aligned}
\hat{x}_{k|k} &= V_{k|k} \left[ C_k^T (D_k D_k^T)^{-1} y_k + (V_{k|k-1})^{-1} \hat{x}_{k|k-1} \right] \\
&= \left[ (B_k B_k^T + A_k V_{k-1|k-1} A_k^T)^{-1} + C_k^T (D_k D_k^T)^{-1} C_k \right]^{-1} \\
&\quad \times \left[ C_k^T (D_k D_k^T)^{-1} y_k + (B_k B_k^T + A_k V_{k-1|k-1})^{-1} A_k \hat{x}_{k-1|k-1} \right] \\
&= \left[ (B_k B_k^T + A_k V_{k-1|k-1} A_k^T)^{-1} + C_k^T (D_k D_k^T)^{-1} C_k \right]^{-1} C_k^T (D_k D_k^T)^{-1} y_k \\
&\quad + \left[ (B_k B_k^T + A_k V_{k-1|k-1} A_k^T)^{-1} + C_k^T (D_k D_k^T)^{-1} C_k \right]^{-1} (B_k B_k^T + A_k V_{k-1|k-1})^{-1} \\
&\quad \times A_k \hat{x}_{k-1|k-1}.
\end{aligned} \tag{B.24}$$

By applying the PosDef Identity (A.3) the following holds.

$$\begin{aligned} & \left[ (B_k B_k^T + A_k V_{k-1|k-1} A_k^T)^{-1} + C_k^T (D_k D_k^T)^{-1} C_k \right]^{-1} C_k^T (D_k D_k^T)^{-1} = \\ & (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T \left[ C_k (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T + (D_k D_k^T) \right]^{-1}. \end{aligned} \quad (\text{B.25})$$

Then, by applying the Woodbury Identity (A.2) the following holds.

$$\begin{aligned} & \left[ (B_k B_k^T + A_k V_{k-1|k-1} A_k^T)^{-1} + C_k^T (D_k D_k^T)^{-1} C_k \right]^{-1} = (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) \\ & - (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T \left[ C_k (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T + (D_k D_k^T) \right]^{-1} \\ & \times C_k (B_k B_k^T + A_k V_{k-1|k-1} A_k^T). \end{aligned} \quad (\text{B.26})$$

Finally, (B.25) and (B.26) are substitute into (B.24) to obtain

$$\begin{aligned} \hat{x}_{k|k} &= (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T \left[ C_k (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T + (D_k D_k^T) \right]^{-1} y_k \\ &+ A \hat{x}_{k-1|k-1} - (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T \\ &\times \left[ C_k (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T + (D_k D_k^T) \right]^{-1} C_k A \hat{x}_{k-1|k-1} \\ &= A \hat{x}_{k-1|k-1} + (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T \\ &\times \left[ C_k (B_k B_k^T + A_k V_{k-1|k-1} A_k^T) C_k^T + (D_k D_k^T) \right]^{-1} (y_k - C_k A \hat{x}_{k-1|k-1}) \end{aligned} \quad (\text{B.27})$$

$$= \hat{x}_{k|k-1} + V_{k|k-1} C_k^T \left[ C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right]^{-1} (y_k - C_k \hat{x}_{k|k-1}). \quad (\text{B.28})$$

Substituting (B.28) into (B.23) yields

$$\begin{aligned} \hat{x}_{k|k-1} &= A_k \hat{x}_{k-1|k-2} + A_k V_{k-1|k-2} C_{k-1}^T \left[ C_{k-1} V_{k-1|k-2} C_{k-1}^T + (D_{k-1}^T D_{k-1}) \right]^{-1} \\ &\times (y_{k-1} - C_{k-1} \hat{x}_{k-1|k-2}). \end{aligned} \quad (\text{B.29})$$

This is the solution of the Kalman Filter.

Next, the term of equation (B.12) which is outside the integral is expanded.

$$\begin{aligned} \frac{|\Sigma_1|^{1/2}}{(2\pi)^{n/2} |D_k| |B_k| |V_{k-1|k-1}|^{1/2}} &= \frac{|\Sigma_1|^{1/2} |V_{k|k}|^{1/2}}{(2\pi)^{n/2} |D_k| |B_k| |V_{k-1|k-1}|^{1/2} |V_{k|k}|^{1/2}} \\ &= \frac{1}{(2\pi)^{n/2} |V_{k|k}|^{1/2}} \exp \left( \log \frac{|\Sigma_1|^{1/2} |V_{k|k}|^{1/2}}{|D_k| |B_k| |V_{k-1|k-1}|^{1/2}} \right). \end{aligned} \quad (\text{B.30})$$

Substituting (B.30) and (B.13) into (B.12) yields

$$\begin{aligned}
\bar{\alpha}_k(x) &= \frac{1}{(2\pi)^{n/2}|V_{k|k}|^{1/2}} \exp\left(- (x - \hat{x}_{k|k})^T \frac{(V_{k|k})^{-1}}{2} (x - \hat{x}_{k|k})\right. \\
&\quad + \hat{x}_{k|k}^T \frac{(V_{k|k})^{-1}}{2} \hat{x}_{k|k} + \hat{x}_{k-1|k-1}^T \left[ (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} \right. \\
&\quad \left. \left. - \frac{(V_{k-1|k-1})^{-1}}{2} \right] \hat{x}_{k-1|k-1} + y_k^T \left[ \frac{1}{2} I_d - \frac{(D_k D_k^T)^{-1}}{2} \right] y_k + \beta_{k-1|k-1} \right. \\
&\quad \left. + \log \frac{|\Sigma_1|^{1/2} |V_{k|k}|^{1/2}}{|D_k| |B_k| |V_{k-1|k-1}|^{1/2}} \right) \\
&= \frac{1}{(2\pi)^{n/2}|V_{k|k}|^{1/2}} \exp\left(- (x - \hat{x}_{k|k})^T \frac{(V_{k|k})^{-1}}{2} (x - \hat{x}_{k|k}) + \beta_{k|k}\right) \quad (\text{B.31})
\end{aligned}$$

where

$$\begin{aligned}
\beta_{k|k} &= \hat{x}_{k|k}^T \frac{(V_{k|k})^{-1}}{2} \hat{x}_{k|k} + \hat{x}_{k-1|k-1}^T \left[ (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} \right. \\
&\quad \left. - \frac{(V_{k-1|k-1})^{-1}}{2} \right] \hat{x}_{k-1|k-1} + y_k^T \left[ \frac{1}{2} I_d - \frac{(D_k D_k^T)^{-1}}{2} \right] y_k \\
&\quad + \log \frac{|\Sigma_1|^{1/2} |V_{k|k}|^{1/2}}{|D_k| |B_k| |V_{k-1|k-1}|^{1/2}} + \beta_{k-1|k-1}. \quad (\text{B.32})
\end{aligned}$$

Now,  $\beta_{k|k}$  will be reshaped in a better form, starting from the expression below,

$$\begin{aligned}
&\hat{x}_{k-1|k-1}^T \left[ (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} - \frac{(V_{k-1|k-1})^{-1}}{2} \right] \hat{x}_{k-1|k-1} \\
&= \hat{x}_{k-1|k-1}^T \left[ (V_{k-1|k-1})^{-1} \frac{\left( (V_{k-1|k-1})^{-1} + A_k^T (B_k B_k^T)^{-1} A_k \right)^{-1}}{2} (V_{k-1|k-1})^{-1} \right. \\
&\quad \left. - \frac{(V_{k-1|k-1})^{-1}}{2} \right] \hat{x}_{k-1|k-1}.
\end{aligned}$$

Next, the Woodbury Identity (A.2) is applied in order to get

$$\begin{aligned}
&= \hat{x}_{k-1|k-1}^T \frac{1}{2} \left[ (V_{k-1|k-1})^{-1} \left( V_{k-1|k-1} - V_{k-1|k-1} A_k^T \left( (B_k B_k^T) + A_k V_{k-1|k-1} A_k^T \right)^{-1} \right. \right. \\
&\quad \left. \left. A_k V_{k-1|k-1} \right) \times (V_{k-1|k-1})^{-1} - (V_{k-1|k-1})^{-1} \right] \hat{x}_{k-1|k-1} \\
&= \hat{x}_{k-1|k-1}^T \frac{1}{2} \left[ (V_{k-1|k-1})^{-1} - A_k^T \left( (B_k B_k^T) + A_k V_{k-1|k-1} A_k^T \right)^{-1} A_k \right. \\
&\quad \left. - (V_{k-1|k-1})^{-1} \right] \hat{x}_{k-1|k-1} \\
&= -\hat{x}_{k-1|k-1}^T \left[ \frac{A_k^T (V_{k|k-1})^{-1} A_k}{2} \right] \hat{x}_{k-1|k-1} \\
&\Rightarrow \hat{x}_{k-1|k-1}^T \left[ (V_{k-1|k-1})^{-1} \frac{\Sigma_1}{2} (V_{k-1|k-1})^{-1} - \frac{(V_{k-1|k-1})^{-1}}{2} \right] \hat{x}_{k-1|k-1} = \\
&\quad -\hat{x}_{k|k-1}^T \left[ \frac{(V_{k|k-1})^{-1}}{2} \right] \hat{x}_{k|k-1}. \quad (\text{B.33})
\end{aligned}$$

Next the expression  $\hat{x}_{k|k}^T \frac{(V_{k|k})^{-1}}{2} \hat{x}_{k|k}$  is reshaped into a better form, by applying the Woodbury Identity (A.2) and the PosDef Identity (A.3).

$$\begin{aligned}
 \hat{x}_{k|k}^T \frac{(V_{k|k})^{-1}}{2} \hat{x}_{k|k} &= \frac{1}{2} \left[ \hat{x}_{k|k-1} + V_{k|k-1} C_k^T \left( C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right)^{-1} y_k \right. \\
 &\quad \left. - V_{k|k-1} C_k^T \left( C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right)^{-1} C_k \hat{x}_{k|k-1} \right]^T \left[ C_k^T (D_k D_k^T)^{-1} C_k + V_{k|k-1}^{-1} \right] \\
 &\quad \left[ \hat{x}_{k|k-1} + V_{k|k-1} C_k^T \left( C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right)^{-1} y_k - V_{k|k-1} C_k^T \left( C_k V_{k|k-1} C_k^T \right. \right. \\
 &\quad \left. \left. + (D_k D_k^T) \right)^{-1} C_k \hat{x}_{k|k-1} \right] \\
 &= \frac{1}{2} \hat{x}_{k|k-1}^T V_{k|k-1}^{-1} \hat{x}_{k|k-1} - \frac{1}{2} y_k^T \left( C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right)^{-1} y_k + \frac{1}{2} y_k^T (D_k D_k^T)^{-1} y_k \\
 &\quad - \frac{1}{2} \hat{x}_{k|k-1}^T C_k^T \left( C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right)^{-1} C_k \hat{x}_{k|k-1} \\
 &\quad + \hat{x}_{k|k-1}^T C_k^T \left( C_k V_{k|k-1} C_k^T + (D_k D_k^T) \right)^{-1} y_k.
 \end{aligned}$$

Finally,

$$\begin{aligned}
 \Rightarrow \hat{x}_{k|k}^T \frac{(V_{k|k})^{-1}}{2} \hat{x}_{k|k} &= \frac{1}{2} \hat{x}_{k|k-1}^T V_{k|k-1}^{-1} \hat{x}_{k|k-1} + \frac{1}{2} y_k^T (D_k D_k^T)^{-1} y_k \\
 &\quad - \frac{1}{2} (y_k - C_k \hat{x}_{k|k-1})^T (C_k V_{k|k-1} C_k^T + D_k D_k^T)^{-1} (y_k - C_k \hat{x}_{k|k-1}).
 \end{aligned} \tag{B.34}$$

Next, the logarithmic term of (B.32) will be reshaped into a better form

$$\begin{aligned}
 &\log \frac{|\Sigma_1|^{1/2} |V_{k|k}|^{1/2}}{|D_k| |B_k| |V_{k-1|k-1}|^{1/2}} \\
 &= \log \frac{|V_{k-1|k-1}^{-1} + A_k^T (B_k B_k^T)^{-1} A_k|^{-1/2} |V_{k|k-1}^{-1} + C_k^T (D_k D_k^T)^{-1} C_k|^{-1/2}}{|D_k| |B_k| |V_{k-1|k-1}|^{1/2}} \\
 &= \log \frac{|V_{k-1|k-1}^{-1} (I + V_{k-1|k-1} A_k^T (B_k B_k^T)^{-1} A_k)|^{-1/2} |V_{k|k-1}^{-1} + C_k^T (D_k D_k^T)^{-1} C_k|^{-1/2}}{|D_k| |B_k| |V_{k-1|k-1}|^{1/2}} \\
 &= \log \frac{|I + V_{k-1|k-1} A_k^T (B_k B_k^T)^{-1} A_k|^{-1/2} |V_{k|k-1}^{-1} + C_k^T (D_k D_k^T)^{-1} C_k|^{-1/2}}{|D_k| |B_k|} \\
 &= \log \frac{|I + (B_k B_k^T)^{-1} A_k V_{k-1|k-1} A_k^T|^{-1/2} |V_{k|k-1}^{-1} + C_k^T (D_k D_k^T)^{-1} C_k|^{-1/2}}{|D_k| |B_k|} \\
 &= \log \frac{|(B_k B_k^T)^{-1} ((B_k B_k^T) + A_k V_{k-1|k-1} A_k^T)|^{-1/2} |V_{k|k-1}^{-1} + C_k^T (D_k D_k^T)^{-1} C_k|^{-1/2}}{|D_k| |B_k|} \\
 &= \log \frac{|(B_k B_k^T)^{-1} V_{k|k-1} (V_{k|k-1}^{-1} + C_k^T (D_k D_k^T)^{-1} C_k)|^{-1/2}}{|D_k| |B_k|}
 \end{aligned}$$

$$\begin{aligned}
&= \log \frac{|(B_k B_k^T)^{-1}(I + V_{k|k-1} C_k^T (D_k D_k^T)^{-1} C_k)|^{-1/2}}{|D_k| |B_k|} \\
&= \log \frac{|(B_k B_k^T)|^{1/2} |(I + (D_k D_k^T)^{-1} C_k V_{k|k-1} C_k^T)|^{-1/2}}{|D_k| |B_k|} \\
&= \log \frac{|(B_k B_k^T)|^{1/2} |(D_k D_k^T)^{-1} (C_k V_{k|k-1} C_k^T + D_k D_k^T)|^{-1/2}}{|D_k| |B_k|} \\
&= \log \frac{|(B_k B_k^T)|^{1/2} |(D_k D_k^T)|^{1/2} |(C_k V_{k|k-1} C_k^T + D_k D_k^T)|^{-1/2}}{|D_k| |B_k|} \\
&= \log |(C_k V_{k|k-1} C_k^T + D_k D_k^T)|^{-1/2} \\
&\Rightarrow \log \frac{|\Sigma_1|^{1/2} |V_{k|k}|^{1/2}}{|D_k| |B_k| |V_{k-1|k-1}|^{1/2}} = -\frac{1}{2} \log |C_k V_{k|k-1} C_k^T + D_k D_k^T|.
\end{aligned}$$

When all the above, (B.33), (B.34) and (B.35), are put together into (B.32) this transforms to

$$\begin{aligned}
\beta_{k|k} &= -\frac{1}{2} \hat{x}_{k|k-1}^T V_{k|k-1}^{-1} \hat{x}_{k|k-1} + \frac{1}{2} \hat{x}_{k|k-1}^T V_{k|k-1}^{-1} \hat{x}_{k|k-1} \\
&\quad + y_k^T \frac{(D_k D_k^T)^{-1}}{2} y_k - \frac{1}{2} (y_k - C_k \hat{x}_{k|k-1})^T (C_k V_{k|k-1} C_k^T + D_k D_k^T)^{-1} \\
&\quad \times (y_k - C_k \hat{x}_{k|k-1}) \\
&\quad + y_k^T \left[ \frac{I_d}{2} - \frac{(D_k D_k^T)^{-1}}{2} \right] y_k - \frac{1}{2} \log |C_k V_{k|k-1} C_k^T + D_k D_k^T| + \beta_{k-1|k-1} \\
&= -\frac{1}{2} (y_k - C_k \hat{x}_{k|k-1})^T (C_k V_{k|k-1} C_k^T + D_k D_k^T)^{-1} (y_k - C_k \hat{x}_{k|k-1}) + \frac{y_k^T y_k}{2} \\
&\quad - \frac{1}{2} \log |C_k V_{k|k-1} C_k^T + D_k D_k^T| + \beta_{k-1|k-1}. \tag{B.35}
\end{aligned}$$

The above expression, (B.35), can also be written as

$$\begin{aligned}
\beta_{k|k} &= -\sum_{i=1}^k \frac{1}{2} (y_i - C_i \hat{x}_{i|i-1})^T (C_i V_{i|i-1} C_i^T + D_i^T D_i)^{-1} (y_i - C_i \hat{x}_{i|i-1}) + \sum_{i=1}^k \frac{y_i^T y_i}{2} \\
&\quad - \frac{1}{2} \sum_{i=1}^k \log |C_i V_{i|i-1} C_i^T + D_i^T D_i|. \tag{B.36}
\end{aligned}$$



## APPENDIX C

### DERIVATIONS OF EXAMPLES OF CHAPTER 3

This Appendix presents the derivations of the solutions of the Examples of Chapter 3.

#### Derivations of Examples of Section 3.3.1 - Estimation of Random Variables

Here, the derivations of the solutions of the examples of Section 3.3.1, are presented.

##### Application 1

The worst case measure,  $\nu^*(y, dx)$ , is given by (3.25). By substituting Bayes' formula  $\eta(y, dx) = \frac{\mu(x, dy)dP_X(x)}{\int_{\mathcal{X}} \mu(x, dy)dP_X(x)}$  into (3.25) the worst case measure becomes

$$\nu^*(y, dx) = \frac{e^{\frac{\ell(x, \Phi(y))}{\hat{s}(y)}} \mu(x, dy)dP_X(x)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\hat{s}(y)}} \mu(x, dy)dP_X(x)}. \quad (\text{C.1})$$

The computation begins from the numerator of (C.1). By substituting (3.55), (3.56) and (3.57) into (C.1) the numerator is given by

$$\begin{aligned} e^{\frac{\ell(x, \Phi(y))}{\hat{s}(y)}} \mu(x, dy)dP_X(x) &= (2\pi)^{-\frac{d}{2}} (2\pi)^{-\frac{n}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \\ &\quad \times \exp \left\{ - (y - Hx)^T \frac{\Sigma_W^{-1}}{2} (y - Hx) \right\} \end{aligned}$$

$$\begin{aligned}
 & -x^T \frac{\Sigma_X^{-1}}{2} x + (x - \Phi(y))^T \frac{U}{\tilde{s}(y)} (x - \Phi(y)) \} dx dy \\
 = & (2\pi)^{-\frac{n}{2}} \exp \left\{ -x^T \left( -\frac{U}{\tilde{s}(y)} + H^T \frac{\Sigma_W^{-1}}{2} H + \frac{\Sigma_X^{-1}}{2} \right) x \right. \\
 & + x^T \left( -2 \frac{U}{\tilde{s}(y)} \Phi(y) + H^T \Sigma_W^{-1} y \right) \\
 & + \Phi(y)^T \frac{U}{\tilde{s}(y)} \Phi(y) - y^T \frac{\Sigma_W^{-1}}{2} y \\
 & \left. + \log \left[ (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \right] \right\} dx dy. \tag{C.2}
 \end{aligned}$$

The function (C.2), is equivalent to

$$\begin{aligned}
 & (2\pi)^{-\frac{n}{2}} \exp \left\{ - (x - m^y(\tilde{s}))^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} (x - m^y(\tilde{s})) + g_1(\tilde{s}, \Phi, y) \right. \\
 & \left. + \log |\Sigma^y(\tilde{s})|^{-\frac{1}{2}} \right\} dx dy \tag{C.3}
 \end{aligned}$$

where

$$\begin{aligned}
 \Sigma^y(\tilde{s}) &= \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{\tilde{s}(y)} \right)^{-1}, \\
 m^y(\tilde{s}) &= \Sigma^y(\tilde{s}) \left( H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi(y) \right), \\
 g_1(\tilde{s}, \Phi, y) &= m^y(\tilde{s})^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} m^y(\tilde{s}) + \Phi(y)^T \frac{U}{\tilde{s}(y)} \Phi(y) - y^T \frac{\Sigma_W^{-1}}{2} y \\
 &+ \log \left[ (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \Sigma^y(\tilde{s})^{\frac{1}{2}} \right].
 \end{aligned}$$

The denominator of (C.1) is equivalent to

$$\begin{aligned}
 & \int_{\mathcal{X}} (2\pi)^{-\frac{n}{2}} \exp \left\{ - (x - m^y(\tilde{s}))^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} (x - m^y(\tilde{s})) + g_1(\tilde{s}, \Phi, y) \right. \\
 & \left. + \log |\Sigma^y(\tilde{s})|^{-\frac{1}{2}} \right\} dx dy. \tag{C.4}
 \end{aligned}$$

Now, given that

$$\int_{\mathcal{X}} \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma^y(\tilde{s})|^{\frac{1}{2}}} e^{-(x - m^y(\tilde{s}))^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} (x - m^y(\tilde{s}))} dx = 1. \tag{C.5}$$

Then, the worst case measure,  $\nu^*(y, dx)$  equals to

$$\nu^*(y, dx) = \frac{\frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma^y(\tilde{s})|^{\frac{1}{2}}} e^{-(x - m^y(\tilde{s}))^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} (x - m^y(\tilde{s})) + g_1(\tilde{s}, \Phi, y)} dx dy}{e^{g_1(\tilde{s}, \Phi, y)} dy} \tag{C.6}$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma^y(\tilde{s})|^{\frac{1}{2}}} e^{-(x - m^y(\tilde{s}))^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} (x - m^y(\tilde{s}))} dx. \tag{C.7}$$



Next, using Bayes' formula the pay-off function (3.27) can be written as

$$L_2(\nu^*, \lambda^*, \tilde{s}^*) = \inf_{\Phi(\cdot)} \inf_{\tilde{s}(\cdot)} \left\{ \tilde{s}(y) \log \frac{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \mu(x, dy) P_X(dx)}{\int_{\mathcal{X}} \mu(x, dy) P_X(dx)} + \tilde{s}(y) R(y) \right\}. \quad (\text{C.8})$$

The function  $\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \mu(x, dy) P_X(dx)$  has already been derived and it is equal to  $e^{g_1(\tilde{s}, \Phi, y)}$ . Next, the following function is derived

$$\begin{aligned} \int_{\mathcal{X}} \mu(x, dy) P_X(dx) &= \int_{\mathcal{X}} (2\pi)^{-\frac{d}{2}} (2\pi)^{-\frac{n}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \\ &\quad \times \exp \left\{ - (y - Hx)^T \frac{\Sigma_W^{-1}}{2} (y - Hx) - x^T \frac{\Sigma_X^{-1}}{2} x \right\} dx dy \end{aligned} \quad (\text{C.9})$$

which is equivalent to

$$\begin{aligned} &\int_{\mathcal{X}} (2\pi)^{-\frac{n}{2}} \exp \left\{ - (x - m_2^y(\tilde{s}))^T \frac{\Sigma_2^y(\tilde{s})^{-1}}{2} (x - m_2^y(\tilde{s})) + g_2(\tilde{s}, y) \right. \\ &\quad \left. + \log |\Sigma_2^y(\tilde{s})|^{-\frac{1}{2}} \right\} dx dy = e^{g_2(\tilde{s}, y)} dy \end{aligned} \quad (\text{C.10})$$

where

$$\begin{aligned} \Sigma_2^y(\tilde{s}) &= (H^T \Sigma_W^{-1} H + \Sigma_X^{-1})^{-1}, \\ m_2^y(\tilde{s}) &= \Sigma_2^y(\tilde{s}) (H^T \Sigma_W^{-1} y), \\ g_2(\tilde{s}, y) &= m_2^y(\tilde{s})^T \frac{\Sigma_2^y(\tilde{s})^{-1}}{2} m_2^y(\tilde{s}) - y^T \frac{\Sigma_W^{-1}}{2} y + \log \left[ (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \Sigma_2^y(\tilde{s})^{\frac{1}{2}} \right]. \end{aligned}$$

Therefore,

$$\frac{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{\tilde{s}(y)}} \mu(x, dy) P_X(dx)}{\int_{\mathcal{X}} \mu(x, dy) P_X(dx)} = \frac{e^{g_1(\tilde{s}, \Phi, y)} dy}{e^{g_2(\tilde{s}, y)} dy} = e^{g_1(\tilde{s}, \Phi, y) - g_2(\tilde{s}, y)}. \quad (\text{C.11})$$

Next, the exponent of the above function is isolated and further manipulated,

$$\begin{aligned} g_1(\tilde{s}, \Phi, y) - g_2(\tilde{s}, y) &= m^y(\tilde{s})^T \frac{\Sigma^y(\tilde{s})^{-1}}{2} m^y(\tilde{s}) + \Phi(y)^T \frac{U}{\tilde{s}(y)} \Phi(y) - y^T \frac{\Sigma_W^{-1}}{2} y \\ &\quad + \log \left[ (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \Sigma^y(\tilde{s})^{\frac{1}{2}} \right] + y^T \frac{\Sigma_W^{-1}}{2} y \\ &\quad - m_2^y(\tilde{s})^T \frac{\Sigma_2^y(\tilde{s})^{-1}}{2} m_2^y(\tilde{s}) - \log \left[ (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \Sigma_2^y(\tilde{s})^{\frac{1}{2}} \right] \\ &= (H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi(y))^T \left[ (H^T \frac{\Sigma_W^{-1}}{2} H + \frac{\Sigma_X^{-1}}{2} - \frac{U}{\tilde{s}(y)})^{-1} \right]^T \end{aligned}$$

$$\begin{aligned}
 & -y^T [\Sigma_W^{-1}]^T H \left[ (H^T \frac{\Sigma_W^{-1}}{2} H + \frac{\Sigma_X^{-1}}{2})^{-1} \right]^T (H^T \Sigma_W^{-1} y) \\
 & + \Phi(y)^T \frac{U}{\tilde{s}(y)} \Phi(y) - \frac{1}{2} \log \frac{|H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{\tilde{s}(y)}|}{|H^T \Sigma_W^{-1} H + \Sigma_X^{-1}|} \\
 = & g(\tilde{s}, \Phi, y) - \frac{1}{2} \log \frac{|\Sigma - \frac{2U}{\tilde{s}(y)}|}{|\Sigma|}
 \end{aligned}$$

where

$$\begin{aligned}
 g(\tilde{s}, \Phi, y) &= (H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi(y))^T \frac{(\Sigma - \frac{2U}{\tilde{s}(y)})^{-1}}{2} (H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi(y)) \\
 &+ \Phi(y)^T \frac{U}{\tilde{s}(y)} \Phi(y) - y^T \Sigma_W^{-1} H \frac{\Sigma^{-1}}{2} H^T \Sigma_W^{-1} y, \\
 \Sigma &= H^T \Sigma_W^{-1} H + \Sigma_X^{-1}.
 \end{aligned}$$

Next, the derived pay-off function has to be differentiated and the derivative set to zero, as shown below

$$\begin{aligned}
 & \frac{d}{d\Phi} L_2(\nu^*, \lambda^*, \tilde{s})|_{\Phi=\Phi^*} = 0, \quad \forall s(\cdot) \\
 & \Rightarrow \frac{d}{d\Phi} \tilde{s}(y) g(\tilde{s}, \Phi, y)|_{\Phi=\Phi^*} = 0.
 \end{aligned}$$

Performing the above differentiation results to

$$\begin{aligned}
 & -2U \left( \Sigma - \frac{2U}{\tilde{s}(y)} \right)^{-1} \left( H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi^*(y) \right) + 2U \Phi^*(y) = 0 \\
 \Rightarrow & \Phi^*(y) = \left( \Sigma - \frac{2U}{\tilde{s}(y)} \right)^{-1} \left( H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi^*(y) \right) \\
 \Rightarrow & \left( \Sigma - \frac{2U}{\tilde{s}(y)} \right) \Phi^*(y) = H^T \Sigma_W^{-1} y - \frac{2U}{\tilde{s}(y)} \Phi^*(y) \\
 \Rightarrow & \Sigma \Phi^*(y) = H^T \Sigma_W^{-1} y \\
 \Rightarrow & \Phi^*(y) = \Sigma^{-1} H^T \Sigma_W^{-1} y = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^T \Sigma_W^{-1} y. \tag{C.12}
 \end{aligned}$$

Finally, applying the PosDef Identity

$$\Phi^*(y) = \Sigma_X H^T \left( H \Sigma_X H^T + \Sigma_W \right)^{-1} y. \tag{C.13}$$

## Application 2

Assuming that the nominal distribution  $\mu(x, dy)$  is given by (3.55), then the true distribution  $\nu(x, dy)$  is given by

$$\nu(x, dy) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_W + \Delta\Sigma_W|^{\frac{1}{2}}} e^{-(y-(H+\Delta H)x)^T \frac{(\Sigma_W + \Delta\Sigma_W)^{-1}}{2} (y-(H+\Delta H)x)} dy. \quad (\text{C.14})$$

The derivation of (3.64) is presented here. The relative entropy between  $\mu(x, dy)$  and  $\nu(x, dy)$  is given by the following expression:

$$\begin{aligned} H(\nu|\mu)(x) &= \int_{\mathcal{Y}} \log \left( \frac{\nu(x, dy)}{\mu(x, dy)} \right) \nu(x, dy) \\ &= \int_{\mathcal{Y}} \nu(x, dy) \log \nu(x, dy) - \int_{\mathcal{Y}} \nu(x, dy) \log \mu(x, dy) \\ &= \int_{\mathcal{Y}} \nu(x, dy) \left\{ -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_W + \Delta\Sigma_W| \right. \\ &\quad \left. - (y - (H + \Delta H)x)^T \frac{(\Sigma_W + \Delta\Sigma_W)^{-1}}{2} (y - (H + \Delta H)x) \right\} \\ &\quad - \int_{\mathcal{Y}} \nu(x, dy) \left\{ -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_W| - (y - Hx)^T \frac{\Sigma_W^{-1}}{2} (y - Hx) \right\} \\ &= \frac{1}{2} \log |\Sigma_W| - \frac{1}{2} \log |\Sigma_W + \Delta\Sigma_W| \\ &\quad - \int_{\mathcal{Y}} \nu(x, dy) (y - (H + \Delta H)x)^T \frac{(\Sigma_W + \Delta\Sigma_W)^{-1}}{2} (y - (H + \Delta H)x) \\ &\quad + \int_{\mathcal{Y}} \nu(x, dy) (y - Hx)^T \frac{\Sigma_W^{-1}}{2} (y - Hx). \end{aligned} \quad (\text{C.15})$$

Next, each term of the above function is treated separately, starting from the term  $\int_{\mathcal{Y}} \nu(x, dy) (y - Hx)^T \frac{\Sigma_W^{-1}}{2} (y - Hx)$  which can be expanded as follows

$$\begin{aligned} &\int_{\mathcal{Y}} \nu(x, dy) (y - Hx)^T \frac{\Sigma_W^{-1}}{2} (y - Hx) = \\ &\int_{\mathcal{Y}} \nu(x, dy) \left( (y - (H + \Delta H)x) - (Hx - (H + \Delta H)x) \right)^T \frac{\Sigma_W^{-1}}{2} \left( (y - (H + \Delta H)x) \right. \\ &\quad \left. - (Hx - (H + \Delta H)x) \right) \\ &= \int_{\mathcal{Y}} \nu(x, dy) \left\{ (y - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2} (y - (H + \Delta H)x) \right. \\ &\quad + (Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2} (Hx - (H + \Delta H)x) \\ &\quad \left. - 2(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2} (y - (H + \Delta H)x) \right\} \end{aligned}$$

$$\begin{aligned}
 &= \int_{\mathcal{Y}} \nu(x, dy)(y - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(y - (H + \Delta H)x) \\
 &+ \int_{\mathcal{Y}} \nu(x, dy)(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(Hx - (H + \Delta H)x) \\
 &- 2 \int_{\mathcal{Y}} \nu(x, dy)(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}y \\
 &+ 2 \int_{\mathcal{Y}} \nu(x, dy)(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(H + \Delta H)x.
 \end{aligned} \tag{C.16}$$

Note that

1.

$$\begin{aligned}
 &\int_{\mathcal{Y}} \nu(x, dy)(y - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(y - (H + \Delta H)x) \\
 &= \text{tr} \left\{ (\Sigma_W + \Delta \Sigma_w) \frac{\Sigma_W^{-1}}{2} \right\}.
 \end{aligned} \tag{C.17}$$

2.

$$\begin{aligned}
 &\int_{\mathcal{Y}} \nu(x, dy)(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(Hx - (H + \Delta H)x) \\
 &= (Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(Hx - (H + \Delta H)x) \\
 &= (\Delta Hx)^T \frac{\Sigma_W^{-1}}{2}(\Delta Hx).
 \end{aligned} \tag{C.18}$$

3.

$$\begin{aligned}
 &-2 \int_{\mathcal{Y}} \nu(x, dy)(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}y \\
 &= -2(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2} \int_{\mathcal{Y}} y \nu(x, dy) \\
 &= -2(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(Hx - (H + \Delta H)x).
 \end{aligned} \tag{C.19}$$

4.

$$\begin{aligned}
 &2 \int_{\mathcal{Y}} \nu(x, dy)(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(H + \Delta H)x \\
 &= 2(Hx - (H + \Delta H)x)^T \frac{\Sigma_W^{-1}}{2}(H + \Delta H)x.
 \end{aligned} \tag{C.20}$$

Therefore, by putting all the above together, (C.16) equals to

$$\text{tr} \left\{ (\Sigma_W + \Delta\Sigma_w) \frac{\Sigma_W^{-1}}{2} \right\} + (\Delta Hx)^T \frac{\Sigma_W^{-1}}{2} (\Delta Hx). \quad (\text{C.21})$$

Similarly,

$$\begin{aligned} & - \int_{\mathcal{Y}} \nu(x, dy) (y - (H + \Delta H)x)^T \frac{(\Sigma_W + \Delta\Sigma_W)^{-1}}{2} (y - (H + \Delta H)x) \\ & = -\text{tr} \left\{ (\Sigma_W + \Delta\Sigma_w) \frac{(\Sigma_W + \Delta\Sigma_W)^{-1}}{2} \right\}. \end{aligned} \quad (\text{C.22})$$

Finally, putting all the above together results to (3.64),

$$\begin{aligned} H(\nu|\mu)(x) & = \int \log \left( \frac{\nu(x, dy)}{\mu(x, dy)} \right) \nu(x, dy) \\ & = \frac{1}{2} \left\{ \log \frac{|\Sigma_W|}{|\Sigma_W + \Delta\Sigma_W|} + \text{tr}((\Sigma_W + \Delta\Sigma_W)(\Sigma_W^{-1} - (\Sigma_W + \Delta\Sigma_W)^{-1})) \right. \\ & \quad \left. + x^T (\Delta H)^T \Sigma_W^{-1} (\Delta H)x \right\}. \end{aligned}$$

### Application 3

The worst case measure,  $dQ_{X,Y}^*(x, y)$ , is given by (3.43). By substituting Bayes' formula  $dP_{X,Y}(y, x) = \mu(x, dy)dP_X(x)$  into (3.43) the worst case measure becomes

$$dQ_{X,Y}^*(x, y) = \frac{e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy) P_X(dx)}{\int_{\mathcal{X} \times \mathcal{Y}} e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy) P_X(dx)}. \quad (\text{C.23})$$

Then the procedure is the same as the one used in Application 1 above, for deriving (3.58), the only difference is that the integration is performed with respect to  $\mathcal{X} \times \mathcal{Y}$  and that  $s$  is not a function of  $y$ . So, the denominator of (C.23) is given by C.2 ( $\tilde{s}(y)$  is replaced by  $s$ ) and is equivalent to

$$\frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma(s)|^{\frac{1}{2}}} e^{-(x - \kappa^y(s))^T \frac{\Sigma(s)^{-1}}{2} (x - \kappa^y(s)) + \theta_1(s, \Phi, y)} dy dx \quad (\text{C.24})$$

where

$$\begin{aligned} \Sigma(s) & = (H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - 2 \frac{U}{s})^{-1}, \\ \kappa^y(s) & = \Sigma(s) (H^T \Sigma_W^{-1} y - \frac{2U}{s} \Phi(y)), \\ \theta_1(s, \Phi, y) & = \kappa^y(s)^T \frac{\Sigma(s)^{-1}}{2} \kappa^y(s) + \Phi(y)^T \frac{U}{s} \Phi(y) - y^T \frac{\Sigma_W^{-1}}{2} y \\ & \quad + \log \left[ (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \right]. \end{aligned}$$

The denominator of (C.23) is given by

$$\int_{\mathcal{X} \times \mathcal{Y}} \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma(s)|^{\frac{1}{2}}} e^{-(x-\kappa^y(s))^T \frac{\Sigma(s)^{-1}}{2} (x-\kappa^y(s)) + \theta_1(s, \Phi, y)} dy dx = \int_{\mathcal{Y}} e^{\theta_1(s, \Phi, y)} dy. \quad (\text{C.25})$$

Therefore, the worst case measure,  $dQ_{X,Y}^*(x, y)$ , is given by

$$\begin{aligned} dQ_{X,Y}^*(x, y) &= \frac{\frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma(s)|^{\frac{1}{2}}} e^{-(x-\kappa^y(s))^T \frac{\Sigma(s)^{-1}}{2} (x-\kappa^y(s)) + \theta_1(s, \Phi, y)} dy dx}{\int_{\mathcal{Y}} e^{\theta_1(s, \Phi, y)} dy} \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma(s)|^{\frac{1}{2}}} e^{-(x-\kappa^y(s))^T \frac{\Sigma(s)^{-1}}{2} (x-\kappa^y(s))} \frac{e^{\theta(s, \Phi, y)}}{\int_{\mathcal{Y}} e^{\theta(s, \Phi, y)} dy} dy dx \quad (\text{C.26}) \end{aligned}$$

where

$$\theta(s, \Phi, y) = \kappa^y(s)^T \frac{\Sigma(s)^{-1}}{2} \kappa^y(s) + \Phi(y)^T \frac{U}{s} \Phi(y) - y^T \frac{\Sigma_W^{-1}}{2} y.$$

Next, the average pay-off is given by (3.42). Using another form of Bayes' formula,

$dP_{X,Y}(y, x) = \eta(y, dx) dP_Y(y)$ , results to,

$$L_4(Q_{X,Y}^*, \lambda^*, s) = s \log \left( \int_{\mathcal{Y}} \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s}} \eta(y, dx) dP_Y(y) \right) + sR. \quad (\text{C.27})$$

The average pay-off (C.27) has to be differentiated and the derivative set to zero. The inner integral will be differentiated as follows

$$\frac{d}{d\Phi} \left[ \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s}} \eta(y, dx) \right] \Big|_{\Phi=\Phi^*} = 0. \quad (\text{C.28})$$

First, the differentiation is performed and then the integration, therefore,

$$\int_{\mathcal{X}} \frac{d}{d\Phi} \left[ e^{\frac{\ell(x, \Phi(y))}{s}} \right] \Big|_{\Phi=\Phi^*} \eta(y, dx) = 0. \quad (\text{C.29})$$

Given (3.57), the above function results to

$$\begin{aligned} \int_{\mathcal{X}} \left( \frac{2U}{s} \Phi^*(y) - \frac{2U}{s} x \right) e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx) &= 0 \\ \Rightarrow \Phi^*(y) &= \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)} = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s}} \mu(x, dy) dP_X(x)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s}} \mu(x, dy) dP_X(x)}. \quad (\text{C.30}) \end{aligned}$$

The function  $e^{\frac{\ell(x, \Phi^*(y))}{s}} \mu(x, dy) dP_X(x)$  has already derived and is equivalent to (C.24) (note that  $\Phi(y) = \Phi^*(y)$ ), therefore

$$\begin{aligned} \Phi^*(y) &= \frac{m^{*,y}(s) e^{\theta_1(s, \Phi^*, y)} dy}{e^{\theta_1(s, \Phi^*, y)} dy} \\ &= \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{s} \right)^{-1} \left( H^T \Sigma_W^{-1} y - \frac{2U}{s} \Phi^*(y) \right). \quad (\text{C.31}) \end{aligned}$$

The above function can be further manipulated, before it reaches its final form. That is,

$$\begin{aligned}
 &\Rightarrow \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{s} \right) \Phi^*(y) = H^T \Sigma_W^{-1} y - \frac{2U}{s} \Phi^*(y) \\
 &\Rightarrow \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} \right) \Phi^*(y) = H^T \Sigma_W^{-1} y \\
 &\Rightarrow \Phi^*(y) = \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^T \Sigma_W^{-1} y.
 \end{aligned} \tag{C.32}$$

Finally, applying the PosDef Identity

$$\Phi^*(y) = \Sigma_X H^T \left( H \Sigma_X H^T + \Sigma_W \right)^{-1} y. \tag{C.33}$$

#### Application 4

The worst case measure,  $\nu^*(y, dx)$ , is given by (3.37). By substituting Bayes' formula  $\eta(y, dx) = \frac{\mu(x, dy) dP_X(x)}{\int_{\mathcal{X}} \mu(x, dy) dP_X(x)}$  into (3.37) the worst case measure becomes

$$\nu^*(y, dx) = \frac{e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)}. \tag{C.34}$$

Then the procedure is the same as the one in Application 1 above. Starting with the numerator, which is given by

$$\begin{aligned}
 e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x) &= (2\pi)^{-\frac{d}{2}} (2\pi)^{-\frac{n}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \\
 &\quad \times \exp \left\{ - (y - Hx)^T \frac{\Sigma_W^{-1}}{2} (y - Hx) - x^T \frac{\Sigma_X^{-1}}{2} x \right. \\
 &\quad \left. + (x - \Phi(y))^T \frac{U}{s} (x - \Phi(y)) \right. \\
 &\quad \left. - (y - \bar{H}x)^T \tilde{U} (y - \bar{H}x) \right\} dx dy \\
 &= (2\pi)^{-\frac{n}{2}} \exp \left\{ - x^T \left( \bar{H}^T \tilde{U} \bar{H} - \frac{U}{s} + H^T \frac{\Sigma_W^{-1}}{2} H \right. \right. \\
 &\quad \left. \left. + \frac{\Sigma_X^{-1}}{2} \right) x + x^T \left( - 2 \frac{U}{s} \Phi(y) + 2 \bar{H}^T \tilde{U} y + H^T \Sigma_W^{-1} y \right) \right. \\
 &\quad \left. + \Phi(y)^T \frac{U}{s} \Phi(y) - y^T \tilde{U} y - y^T \frac{\Sigma_W^{-1}}{2} y \right. \\
 &\quad \left. + \log \left[ (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{-\frac{1}{2}} \right] \right\} dx dy.
 \end{aligned} \tag{C.35}$$

The above function (C.35) is equivalent with the following function

$$(2\pi)^{-\frac{n}{2}} \exp \left\{ - (x - \tilde{m}^y(s))^T \frac{\tilde{\Sigma}(s)^{-1}}{2} (x - \tilde{m}^y(s)) + g_1(s, \Phi, y) \right\}$$

$$+ \log |\tilde{\Sigma}(s)|^{-\frac{1}{2}} \} dx dy \quad (\text{C.36})$$

where

$$\begin{aligned} \tilde{\Sigma}(s) &= \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\bar{H}^T \tilde{U} \bar{H} - \frac{2U}{s} \right)^{-1}, \\ \tilde{m}^y(s) &= \tilde{\Sigma}(s) \left( H^T \Sigma_W^{-1} y + 2\bar{H}^T \tilde{U} y - \frac{2U}{s} \Phi(y) \right). \end{aligned}$$

The denominator of (C.34) is equivalent to

$$\begin{aligned} \int_{\mathcal{X}} (2\pi)^{-\frac{n}{2}} \exp \left\{ - (x - \tilde{m}^y(s))^T \frac{\tilde{\Sigma}(s)^{-1}}{2} (x - \tilde{m}^y(s)) + g_1(s, \Phi, y) \right. \\ \left. + \log |\tilde{\Sigma}(s)|^{-\frac{1}{2}} \right\} dx dy = e^{g_1(s, \Phi, y)} dy. \end{aligned} \quad (\text{C.37})$$

Therefore, the worst case measure, (C.34), is finally given by

$$\begin{aligned} \nu^*(y, dx) &= \frac{\frac{1}{(2\pi)^{\frac{n}{2}} |\tilde{\Sigma}(s)|^{\frac{1}{2}}} e^{-(x - \tilde{m}^y(s))^T \frac{\tilde{\Sigma}(s)^{-1}}{2} (x - \tilde{m}^y(s)) + g_1(s, \Phi, y)} dx dy}{e^{g_1(s, \Phi, y)} dy} \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\tilde{\Sigma}(s)|^{\frac{1}{2}}} e^{-(x - \tilde{m}^y(s))^T \frac{\tilde{\Sigma}(s)^{-1}}{2} (x - \tilde{m}^y(s))} dx. \end{aligned} \quad (\text{C.38})$$

Next, the average pay-off,  $L_3^R(\nu^*, \lambda^*, s^*)$  is given by (3.39)

$$L_3^R(\nu^*, \lambda^*, s^*) = \inf_{s \geq 0} \int_{\mathcal{Y}} s \log \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \eta(y, dx) dP_Y(y) + s\bar{R}.$$

The average pay-off (3.39) has to be differentiated and the derivative set to zero.

Similar to Application 3 above, the inner integral will be differentiated as follows

$$\frac{d}{d\Phi} \left[ \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \eta(y, dx) \right] \Big|_{\Phi = \Phi^*} = 0. \quad (\text{C.39})$$

First, the differentiation is performed and then the integration, therefore,

$$\int_{\mathcal{X}} \frac{d}{d\Phi} \left[ e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \right] \Big|_{\Phi = \Phi^*} \eta(y, dx) = 0. \quad (\text{C.40})$$

Given (3.57), the above function results in

$$\int_{\mathcal{X}} \left( \frac{2U}{s} \Phi^*(y) - \frac{2U}{s} x \right) e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx) = 0$$



$$\begin{aligned}
 \Rightarrow \Phi^*(y) &= \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx)} \\
 &= \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)}. \tag{C.41}
 \end{aligned}$$

The function  $e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)$  is derived and is equivalent to (C.36) (note that  $\Phi(y) = \Phi^*(y)$ ). Therefore

$$\begin{aligned}
 \Phi^*(y) &= \frac{\tilde{m}^{*,y}(s) e^{g_1(s, \Phi^*, y)} dy}{e^{g_1(s, \Phi^*, y)} dy} \\
 &= \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\bar{H}^T \tilde{U} \bar{H} - \frac{2U}{s} \right)^{-1} \left( H^T \Sigma_W^{-1} y + 2\bar{H}^T \tilde{U} y - \frac{2U}{s} \Phi^*(y) \right). \tag{C.42}
 \end{aligned}$$

The above function can be further manipulated, before it reaches its final form. That is,

$$\begin{aligned}
 \Rightarrow \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\bar{H}^T \tilde{U} \bar{H} - \frac{2U}{s} \right) \Phi^*(y) &= H^T \Sigma_W^{-1} y + 2\bar{H}^T \tilde{U} y - \frac{2U}{s} \Phi^*(y) \\
 \Rightarrow \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\bar{H}^T \tilde{U} \bar{H} \right) \Phi^*(y) &= H^T \Sigma_W^{-1} y + 2\bar{H}^T \tilde{U} y \\
 \Rightarrow \Phi^*(y) &= \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} + 2\bar{H}^T \tilde{U} \bar{H} \right)^{-1} \left( H^T \Sigma_W^{-1} y + 2\bar{H}^T \tilde{U} y \right). \tag{C.43}
 \end{aligned}$$

## Derivation of Example of Section 3.3.2 - Estimation of a Sequence of Random Variables

Just like in Application 1 of Section 3.3.1, the worst case measure,  $\nu^*(y, dx)$ , is given by (C.1). Similarly, the computation begins from the numerator of (C.1). By substituting (3.78), (3.79) and (3.80) into (C.1) the numerator is given by

$$\begin{aligned}
 e^{\frac{\ell(x, \Phi(y^m))}{\tilde{s}(y)}} \mu(x, dy^m) dP_X(x) &= \left( (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} \right)^{m+1} (2\pi)^{-\frac{n}{2}} |\Sigma_X|^{-\frac{1}{2}} \\
 &\quad \times \exp \left\{ - \sum_{i=0}^m (y_i - Hx)^T \frac{\Sigma_W^{-1}}{2} (y_i - Hx) - x^T \frac{\Sigma_X^{-1}}{2} x \right. \\
 &\quad \left. + \sum_{i=0}^{m-1} (x - \Phi_i^*(y^i))^T \frac{U}{\tilde{s}(y)} (x - \Phi_i^*(y^i)) \right\}
 \end{aligned}$$

$$\begin{aligned}
 & +(x - \Phi_m(y^m))^T \frac{U}{\tilde{s}(y)} (x - \Phi_m(y^m)) \Big\} dx dy^m \\
 = & (2\pi)^{-\frac{n}{2}} \exp \left\{ -x^T \left( (m+1)H^T \frac{\Sigma_W^{-1}}{2} H + \frac{\Sigma_X^{-1}}{2} \right. \right. \\
 & \left. \left. - (m+1) \frac{U}{\tilde{s}(y)} \right) x + x^T \left( -2 \frac{U}{\tilde{s}(y)} \sum_{i=0}^{(m-1)} \Phi_i^*(y^i) \right. \right. \\
 & \left. \left. - \frac{2U\Phi_m(y^m)}{\tilde{s}(y)} + H^T \Sigma_W^{-1} \sum_{i=0}^m y_i \right) - \sum_{i=1}^N y_i^T \frac{\Sigma_W^{-1}}{2} y_i \right. \\
 & \left. + \sum_{i=0}^{(m-1)} \Phi_i^*(y^i)^T \frac{U}{\tilde{s}(y)} \Phi_i^*(y^i) + \Phi_m(y^m)^T \frac{U}{\tilde{s}(y)} \Phi_m(y^m) \right. \\
 & \left. + \log \left[ (2\pi)^{-\frac{d(m+1)}{2}} |\Sigma_W|^{-\frac{m+1}{2}} |\Sigma_X|^{-\frac{1}{2}} \right] \right\} dx dy^m. \quad (\text{C.44})
 \end{aligned}$$

The above function (C.44) is equivalent with the following function

$$\begin{aligned}
 & (2\pi)^{-\frac{n}{2}} \exp \left\{ - (x - \hat{m}^y(\tilde{s}))^T \frac{\hat{\Sigma}^y(\tilde{s})^{-1}}{2} (x - \hat{m}^y(\tilde{s})) + \tilde{g}_1(\tilde{s}, \Phi, y) \right. \\
 & \left. + \log |\hat{\Sigma}^y(\tilde{s})|^{-\frac{1}{2}} \right\} dx dy \quad (\text{C.45})
 \end{aligned}$$

where

$$\begin{aligned}
 \hat{\Sigma}^y(\tilde{s}) &= \left( NH^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2(m+1)U}{\tilde{s}(y)} \right)^{-1}, \\
 \hat{m}^y(\tilde{s}) &= \hat{\Sigma}^y(\tilde{s}) \left( H^T \Sigma_W^{-1} \sum_{i=0}^m y_i - \frac{2U}{\tilde{s}(y)} \sum_{i=0}^{m-1} \Phi_i^*(y^i) - \frac{2U\Phi_m(y^m)}{\tilde{s}(y)} \right).
 \end{aligned}$$

The denominator of (C.1) is equivalent to

$$\begin{aligned}
 & \int_{\mathcal{X}} (2\pi)^{-\frac{n}{2}} \exp \left\{ - (x - \hat{m}^y(\tilde{s}))^T \frac{\hat{\Sigma}^y(\tilde{s})^{-1}}{2} (x - \hat{m}^y(\tilde{s})) + \tilde{g}_1(\tilde{s}, \Phi, y) \right. \\
 & \left. + \log |\hat{\Sigma}^y(\tilde{s})|^{-\frac{1}{2}} \right\} dx dy^m = e^{\tilde{g}_1(\tilde{s}, \Phi, y)} dy^m. \quad (\text{C.46})
 \end{aligned}$$

Therefore, the worst case measure, (C.1), is finally given by

$$\begin{aligned}
 \nu^*(y, dx) &= \frac{\frac{1}{(2\pi)^{\frac{n}{2}} |\hat{\Sigma}^y(\tilde{s})|^{-\frac{1}{2}}} e^{-(x - \hat{m}^y(\tilde{s}))^T \frac{\hat{\Sigma}^y(\tilde{s})^{-1}}{2} (x - \hat{m}^y(\tilde{s})) + \tilde{g}_1(\tilde{s}, \Phi, y)} dx dy}{e^{\tilde{g}_1(\tilde{s}, \Phi, y)} dy} \\
 &= \frac{1}{(2\pi)^{\frac{n}{2}} |\hat{\Sigma}^y(\tilde{s})|^{-\frac{1}{2}}} e^{-(x - \hat{m}^y(\tilde{s}))^T \frac{\hat{\Sigma}^y(\tilde{s})^{-1}}{2} (x - \hat{m}^y(\tilde{s}))}. \quad (\text{C.47})
 \end{aligned}$$

Next, the average pay-off function  $L_2(\nu^*, \lambda^*, s)$  given by (3.27) has to be differentiated and the derivative set to zero. The inner integral will be differentiated as follows

$$\frac{d}{d\Phi_m(y^m)} \left[ \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y^m))}{\tilde{s}(y)}} \eta(y^m, dx) \right] \Big|_{\Phi_m(y^m) = \Phi_m^*(y^m)} = 0. \quad (\text{C.48})$$

First, the differentiation is performed and then the integration, therefore,

$$\int_{\mathcal{X}} \frac{d}{d\Phi_m(y^m)} \left[ e^{\frac{\ell(x, \Phi(y^m))}{\tilde{s}(y)}} \right] \Big|_{\Phi_m(y^m) = \Phi_m^*(y^m)} \eta(y^m, dx) = 0. \quad (\text{C.49})$$

Given (3.80), the above function results to

$$\begin{aligned} & \int_{\mathcal{X}} \left( \frac{2U}{\tilde{s}(y)} \Phi_m^*(y^m) - \frac{2U}{\tilde{s}(y)} x \right) e^{\frac{\ell(x, \Phi^*(y^m))}{\tilde{s}(y)}} \eta(y^m, dx) = 0 \\ \Rightarrow \Phi_m^*(y^m) &= \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y^m))}{\tilde{s}(y)}} \eta(y^m, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y^m))}{\tilde{s}(y)}} \eta(y^m, dx)} = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y^m))}{\tilde{s}(y)}} \mu(x, dy^m) dP_X(x)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y^m))}{\tilde{s}(y)}} \mu(x, dy^m) dP_X(x)}. \end{aligned} \quad (\text{C.50})$$

The function  $e^{\frac{\ell(x, \Phi^*(y^m))}{\tilde{s}(y)}} \mu(x, dy^m) dP_X(x)$  has already derived and is equivalent to (C.45) (note that  $\Phi_m(y^m) = \Phi_m^*(y^m)$ ), therefore

$$\begin{aligned} \Phi_m^*(y^m) &= \frac{\hat{m}^{*,y}(\tilde{s}) e^{\tilde{g}_1(s, \Phi^*, y)} dy^m}{e^{\tilde{g}_1(s, \Phi^*, y)} dy^m} \\ &= \left( (m+1) H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2(m+1)U}{\tilde{s}(y)} \right)^{-1} \\ &\quad \times \left( H^T \Sigma_W^{-1} \sum_{i=0}^m y_i - \frac{2U}{\tilde{s}(y)} \sum_{i=0}^{m-1} \Phi_i^*(y^i) - \frac{2U}{\tilde{s}(y)} \Phi_m^*(y^m) \right). \end{aligned} \quad (\text{C.51})$$

The above function can be further manipulated, before it reaches its final form. That is,

$$\begin{aligned} \Rightarrow & \left( (m+1) H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2(m+1)U}{\tilde{s}(y)} \right) \Phi_m^*(y^m) \\ &= H^T \Sigma_W^{-1} \sum_{i=0}^m y_i - \frac{2U}{\tilde{s}(y)} \sum_{i=0}^{m-1} \Phi_i^*(y^i) - \frac{2U}{\tilde{s}(y)} \Phi_m^*(y^m) \\ \Rightarrow \Phi_m^*(y^m) &= \left( (m+1) H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2mU}{\tilde{s}(y)} \right)^{-1} \left( H^T \Sigma_W^{-1} \sum_{i=0}^m y_i \right. \\ &\quad \left. - \frac{2U}{\tilde{s}(y)} \sum_{i=0}^{m-1} \Phi_i^*(y^i) \right). \end{aligned} \quad (\text{C.52})$$

## Derivations of Examples of Section 3.4.2 - Estimation from MIMO communication Systems

Here, the derivation of the results of the examples of Section 3.4.2, are presented.

### Application 1

The procedure followed here is very similar to the one applied in Application 4 of Section 3.3.1. The worst case measure,  $\nu^*(y, dx)$ , is given by (C.34). By substituting (3.87), (3.88), (3.89) and (3.90) into (C.1) the numerator is given by

$$\begin{aligned}
 e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x) &= (2\pi)^{-d} (2\pi)^{-n} |\Sigma_W|^{-1} |\Sigma_X|^{-1} \\
 &\times \exp \left\{ - (y - Hx)^\dagger \Sigma_W^{-1} (y - Hx) - x^\dagger \Sigma_X^{-1} x \right. \\
 &+ (x - \Phi(y))^\dagger \frac{U}{s} (x - \Phi(y)) \\
 &\left. - (y - \bar{H}x)^\dagger \tilde{U} (y - \bar{H}x) \right\} dx dy \\
 &= (2\pi)^{-n} \exp \left\{ - x^\dagger \left( \bar{H}^\dagger \tilde{U} \bar{H} - \frac{U}{s} + H^T \Sigma_W^{-1} H \right. \right. \\
 &\left. \left. + \Sigma_X^{-1} \right) x + x^\dagger \left( - 2\bar{H}^\dagger \frac{U}{s} \Phi(y) + 2\tilde{U} y + 2H^T \Sigma_W^{-1} y \right) \right. \\
 &\left. + \Phi(y)^\dagger \frac{U}{s} \Phi(y) - y^\dagger \tilde{U} y - y^\dagger \Sigma_W^{-1} y \right. \\
 &\left. + \log \left[ (2\pi)^{-d} |\Sigma_W|^{-1} |\Sigma_X|^{-1} \right] \right\} dx dy. \tag{C.53}
 \end{aligned}$$

The above function (C.53) is equivalent to the following function

$$\begin{aligned}
 (2\pi)^{-n} \exp \left\{ - (x - \tilde{m}_M(s))^\dagger \tilde{\Sigma}_M(s)^{-1} (x - \tilde{m}_M(s)) + q(s, \Phi, y) \right. \\
 \left. + \log |\tilde{\Sigma}_M(s)|^{-1} \right\} dx dy \tag{C.54}
 \end{aligned}$$

where

$$\begin{aligned}
 \tilde{\Sigma}_M(s) &= \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \bar{H}^\dagger \tilde{U} \bar{H} - \frac{U}{s} \right)^{-1}, \\
 \tilde{m}_M(s) &= \tilde{\Sigma}(s) \left( H^\dagger \Sigma_W^{-1} y + \bar{H}^\dagger \tilde{U} y - \frac{U}{s} \Phi(y) \right).
 \end{aligned}$$

The denominator of (C.34) is equivalent to

$$\int_{\mathcal{X}} (2\pi)^{-n} \exp \left\{ - (x - \tilde{m}_M(s))^\dagger \tilde{\Sigma}_M(s)^{-1} (x - \tilde{m}_M(s)) + q(s, \Phi, y) \right.$$

$$+ \log |\tilde{\Sigma}_M(s)|^{-1} \} dx dy = e^{q(s, \Phi, y)} dy. \quad (\text{C.55})$$

Therefore, the worst case measure, (C.34), is finally given by

$$\begin{aligned} \nu^*(y, dx) &= \frac{(2\pi)^{-n} |\tilde{\Sigma}_M(s)|^{-1} e^{-(x - \tilde{m}_M(s))^\dagger \tilde{\Sigma}_M(s)^{-1} (x - \tilde{m}_M(s)) + q(s, \Phi, y)} dx dy}{e^{q(s, \Phi, y)} dy} \\ &= (2\pi)^{-n} |\tilde{\Sigma}_M(s)|^{-1} e^{-(x - \tilde{m}_M(s))^\dagger \tilde{\Sigma}_M(s)^{-1} (x - \tilde{m}_M(s))} dx. \end{aligned} \quad (\text{C.56})$$

Next, the average pay-off,  $L_3^R(\nu^*, \lambda^*, s^*)$  is given by (3.39)

$$L_3^R(\nu^*, \lambda^*, s^*) = \inf_{s \geq 0} \int_{\mathcal{Y}} s \log \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \eta(y, dx) dP_Y(y) + s\bar{R}.$$

The average pay-off (3.39) has to be differentiated and the derivative set to zero. Similar to Application 4 of Section 3.3.1, the inner integral will be differentiated as follows

$$\frac{d}{d\Phi} \left[ \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \eta(y, dx) \right] \Big|_{\Phi = \Phi^*} = 0. \quad (\text{C.57})$$

First, the differentiation is performed and then the integration, therefore,

$$\int_{\mathcal{X}} \frac{d}{d\Phi} \left[ e^{\frac{\ell(x, \Phi(y))}{s} - R(x, y)} \right] \Big|_{\Phi = \Phi^*} \eta(y, dx) = 0. \quad (\text{C.58})$$

Given (3.89), the above function is

$$\begin{aligned} \int_{\mathcal{X}} \left( \frac{2U}{s} \Phi^*(y) - \frac{2U}{s} x \right) e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx) &= 0 \\ \Rightarrow \Phi^*(y) &= \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \eta(y, dx)} \\ &= \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)}. \end{aligned} \quad (\text{C.59})$$

The function  $e^{\frac{\ell(x, \Phi^*(y))}{s} - R(x, y)} \mu(x, dy) dP_X(x)$  has already derived and is equivalent to (C.54) (note that  $\Phi(y) = \Phi^*(y)$ ), therefore

$$\begin{aligned} \Phi^*(y) &= \frac{\tilde{m}_M^*(s) e^{q(s, \Phi^*, y)} dy}{e^{q(s, \Phi^*, y)} dy} \\ &= \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \bar{H}^\dagger \tilde{U} \bar{H} - \frac{U}{s} \right)^{-1} \left( H^\dagger \Sigma_W^{-1} y + \bar{H}^\dagger \tilde{U} y - \frac{U}{s} \Phi^*(y) \right). \end{aligned} \quad (\text{C.60})$$

The above function can be further manipulated, before it reaches its final form. That is,

$$\begin{aligned}
 &\Rightarrow \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \bar{H}^\dagger \tilde{U} \bar{H} - \frac{U}{s} \right) \Phi^*(y) = H^\dagger \Sigma_W^{-1} y + \bar{H}^\dagger \tilde{U} y - \frac{U}{s} \Phi^*(y) \\
 &\Rightarrow \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \bar{H}^\dagger \tilde{U} \bar{H} \right) \Phi^*(y) = H^\dagger \Sigma_W^{-1} y + \bar{H}^\dagger \tilde{U} y \\
 &\Rightarrow \Phi^*(y) = \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} + \bar{H}^\dagger \tilde{U} \bar{H} \right)^{-1} \left( H^\dagger \Sigma_W^{-1} y + \bar{H}^\dagger \tilde{U} y \right). \tag{C.61}
 \end{aligned}$$

## Application 2

The procedure followed here is very similar to the one applied in Application 3 of Section 3.3.1. The worst case measure,  $dQ_{X,Y}^*(x, y)$ , is given by (C.23). By substituting (3.87), (3.88), (3.89) and (3.90) into (C.23) the numerator is given by

$$\begin{aligned}
 e^{\frac{\ell(x, \Phi(y))}{s}} \mu(x, dy) dP_X(x) &= (2\pi)^{-d} (2\pi)^{-n} |\Sigma_W|^{-1} |\Sigma_X|^{-1} \\
 &\quad \times \exp \left\{ - (y - Hx)^\dagger \Sigma_W^{-1} (y - Hx) - x^\dagger \Sigma_X^{-1} x \right. \\
 &\quad \left. + (x - \Phi(y))^\dagger \frac{U}{s} (x - \Phi(y)) \right\} dx dy \\
 &= (2\pi)^{-n} \exp \left\{ - x^\dagger \left( H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{U}{s} \right) x \right. \\
 &\quad \left. + x^\dagger \left( - 2 \frac{U}{s} \Phi(y) + 2 H^T \Sigma_W^{-1} y \right) \right. \\
 &\quad \left. + \Phi(y)^\dagger \frac{U}{s} \Phi(y) - y^\dagger \Sigma_W^{-1} y \right. \\
 &\quad \left. + \log \left[ (2\pi)^{-d} |\Sigma_W|^{-1} |\Sigma_X|^{-1} \right] \right\} dx dy. \tag{C.62}
 \end{aligned}$$

The above function (C.62) is equivalent with the following function

$$\begin{aligned}
 &(2\pi)^{-n} \exp \left\{ - (x - \kappa_M^y(s))^\dagger \Sigma_M(s)^{-1} (x - \kappa_M^y(s)) + \tilde{\theta}_1(s, \Phi, y) \right. \\
 &\quad \left. + \log |\Sigma_M(s)|^{-1} \right\} dx dy \tag{C.63}
 \end{aligned}$$

where

$$\begin{aligned}
 \Sigma_M(s) &= \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{U}{s} \right)^{-1}, \\
 \kappa_M^y(s) &= \Sigma_M(s) \left( H^\dagger \Sigma_W^{-1} y - \frac{U}{s} \Phi(y) \right), \\
 \tilde{\theta}_1(s, \Phi, y) &= \kappa_M^y(s)^\dagger \Sigma_M(s)^{-1} \kappa_M^y(s) + \Phi(y)^\dagger \frac{U}{s} \Phi(y) - y^\dagger \Sigma_W^{-1} y \\
 &\quad + \log \left[ (2\pi)^{-d} |\Sigma_W|^{-1} |\Sigma_X|^{-1} \right].
 \end{aligned}$$

The denominator of (C.23) is equivalent to

$$\begin{aligned} & \int_{\mathcal{X} \times \mathcal{Y}} (2\pi)^{-n} |\Sigma_M(s)|^{-1} e^{-(x - \kappa_M^y(s))^\dagger \Sigma_M(s)^{-1} (x - \kappa_M^y(s)) + \tilde{\theta}_1(s, \Phi, y)} dx dy \\ &= \int_{\mathcal{Y}} e^{\tilde{\theta}_1(s, \Phi, y)} dy. \end{aligned} \quad (\text{C.64})$$

Therefore, the worst case measure,  $dQ_{X,Y}^*(x, y)$ , is given by

$$\begin{aligned} dQ_{X,Y}^*(x, y) &= \frac{(2\pi)^{-n} |\Sigma_M(s)|^{-1} e^{-(x - \kappa_M^y(s))^\dagger \Sigma_M(s)^{-1} (x - \kappa_M^y(s)) + \tilde{\theta}_1(s, \Phi, y)} dx dy}{\int_{\mathcal{Y}} e^{\tilde{\theta}_1(s, \Phi, y)} dy} \\ &= (2\pi)^{-n} |\Sigma_M(s)|^{-1} e^{-(x - \kappa_M^y(s))^\dagger \Sigma_M(s)^{-1} (x - \kappa_M^y(s))} \frac{e^{\tilde{\theta}(s, \Phi, y)}}{\int_{\mathcal{Y}} e^{\tilde{\theta}(s, \Phi, y)} dy} dy dx \end{aligned} \quad (\text{C.65})$$

where

$$\tilde{\theta}(s, \Phi, y) = \kappa_M^y(s)^\dagger \Sigma_M(s)^{-1} \kappa_M^y(s) + \Phi(y)^\dagger \frac{U}{s} \Phi(y) - y^\dagger \Sigma_W^{-1} y.$$

Next, the average pay-off, given by (C.27) has to be differentiated and the derivative set to zero. The inner integral will be differentiated as follows.

$$\frac{d}{d\Phi} \left[ \int_{\mathcal{X}} e^{\frac{\ell(x, \Phi(y))}{s}} \eta(y, dx) \right] \Big|_{\Phi = \Phi^*} = 0. \quad (\text{C.66})$$

First, the differentiation is performed and then the integration, therefore,

$$\int_{\mathcal{X}} \frac{d}{d\Phi} \left[ e^{\frac{\ell(x, \Phi(y))}{s}} \right] \Big|_{\Phi = \Phi^*} \eta(y, dx) = 0. \quad (\text{C.67})$$

Given (3.89), the above function results to

$$\begin{aligned} & \int_{\mathcal{X}} \left( \frac{2U}{s} \Phi^*(y) - \frac{2U}{s} x \right) e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx) = 0 \\ \Rightarrow \Phi^*(y) &= \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s}} \eta(y, dx)} = \frac{\int_{\mathcal{X}} x e^{\frac{\ell(x, \Phi^*(y))}{s}} \mu(x, dy) dP_X(x)}{\int_{\mathcal{X}} e^{\frac{\ell(x, \Phi^*(y))}{s}} \mu(x, dy) dP_X(x)}. \end{aligned} \quad (\text{C.68})$$

The function  $e^{\frac{\ell(x, \Phi^*(y))}{s}} \mu(x, dy) dP_X(x)$  has already derived and is equivalent to (C.63) (note that  $\Phi(y) = \Phi^*(y)$ ). Therefore

$$\begin{aligned} \Phi^*(y) &= \frac{\kappa_M^{*,y}(s) e^{\tilde{\theta}_1(s, \Phi^*, y)} dy}{e^{\tilde{\theta}_1(s, \Phi^*, y)} dy} \\ &= \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{U}{s} \right)^{-1} \left( H^\dagger \Sigma_W^{-1} y - \frac{U}{s} \Phi^*(y) \right) \end{aligned} \quad (\text{C.69})$$

The above function can be further manipulated, before it reaches its final form. That is,

$$\begin{aligned}
&\Rightarrow \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{U}{s} \right) \Phi^*(y) = H^\dagger \Sigma_W^{-1} y - \frac{U}{s} \Phi^*(y) \\
&\Rightarrow \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} \right) \Phi^*(y) = H^\dagger \Sigma_W^{-1} y \\
&\Rightarrow \Phi^*(y) = \left( H^\dagger \Sigma_W^{-1} H + \Sigma_X^{-1} \right)^{-1} H^\dagger \Sigma_W^{-1} y. \tag{C.70}
\end{aligned}$$

Finally, applying the PosDef Identity

$$\Phi^*(y) = \Sigma_X H^\dagger \left( H \Sigma_X H^\dagger + \Sigma_W \right)^{-1} y \tag{C.71}$$



## APPENDIX D

### PROOF OF THEOREM 4.4.4

First the following parameters are defined,  $\bar{\sigma}^{2,m} = \frac{\sigma^2}{1+2\sigma^2 K_i^m}$  and  $\epsilon^m = V_i(y^m)\bar{\sigma}^{2,m}$  which are used throughout the derivation.

(a) When (4.63) is substituted into (4.64), the following is obtained.

$$\begin{aligned}
\hat{\Lambda}(t_m) &= \prod_{i=1}^N \left[ \int_0^\infty \frac{r_i}{\sigma^2} \exp\left(-\frac{r_i^2}{2\sigma^2}\right) \exp(-r_i^2 K_i^m) \right. \\
&\quad \times \left. \frac{1}{2\pi} \int_0^{2\pi} \exp\left(r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m))\right) d\theta_i dr_i \right] \\
&\quad \times \exp\left(\frac{1}{2} \sum_{k=0}^m y^2(t_k)[1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)|\right) \\
&= \prod_{i=1}^N \left[ \int_0^\infty \frac{r_i}{\sigma^2} \exp\left(-\frac{r_i^2}{2\left(\frac{\sigma^2}{1+2\sigma^2 K_i^m}\right)}\right) I_0\left(r_i V_i(y^m)\right) dr_i \right] \\
&\quad \times \exp\left(\frac{1}{2} \sum_{k=0}^m y^2(t_k)[1 - D^{-2}(t_k)] - \sum_{k=0}^m \log |D(t_k)|\right) \tag{D.1}
\end{aligned}$$

where  $I_0(\cdot)$  is the modified bessel function of the first kind and zeroth order defined by  $I_0(x) \triangleq \frac{1}{2\pi} \int_{-\pi}^{\pi} \exp(x \cos \alpha) d\alpha$ .

Next, using  $\bar{\sigma}^{2,m}$  and  $\epsilon^m$  the integral  $\int_0^\infty \frac{r_i}{\sigma^2} \exp\left(-\frac{r_i^2}{2\left(\frac{\sigma^2}{1+2\sigma^2 K_i^m}\right)}\right) I_0\left(r_i V_i(y^m)\right) dr_i$  is computed as follows

$$\begin{aligned}
&\int_0^\infty \frac{r_i}{\sigma^2} \exp\left(-\frac{r_i^2}{2\left(\frac{\sigma^2}{1+2\sigma^2 K_i^m}\right)}\right) I_0\left(r_i V_i(y^m)\right) dr_i = \int_0^\infty \frac{r_i}{\sigma^2} \exp\left(-\frac{r_i^2}{2\bar{\sigma}^{2,m}}\right) I_0\left(\frac{r_i \epsilon^m}{\bar{\sigma}^{2,m}}\right) dr_i \\
&= \frac{\bar{\sigma}^{2,m}}{\sigma^2} \exp\left(\frac{\epsilon^{2,m}}{2\bar{\sigma}^{2,m}}\right) \int_0^\infty \frac{r_i}{\bar{\sigma}^{2,m}} \exp\left(-\frac{r_i^2}{2\bar{\sigma}^{2,m}}\right) I_0\left(\frac{r_i \epsilon^m}{\bar{\sigma}^{2,m}}\right) dr_i = \frac{\bar{\sigma}^{2,m}}{\sigma^2} \exp\left(\frac{\epsilon^{2,m}}{2\bar{\sigma}^{2,m}}\right)
\end{aligned}$$

$$= \frac{1}{1 + 2\sigma^2 K_i^m} \exp\left(\frac{V_i^2(y^m)\sigma^2}{2 + 4\sigma^2 K_i^m}\right). \quad (\text{D.2})$$

Note that in (D.2) the fact that  $\frac{r_i}{\bar{\sigma}^{2,m}} \exp\left(-\frac{(r_i^2 + \epsilon^{2,m})}{2\bar{\sigma}^{2,m}}\right) I_0\left(\frac{r_i \epsilon^m}{\bar{\sigma}^{2,m}}\right)$  is a Rice probability distribution function which integrates to 1, is being used. By Substituting (D.2) into (D.1), equation (4.65) is obtained.

(b) The normalized conditional density  $p_N(t_m, \boldsymbol{\theta}, \mathbf{r} | \mathcal{Y}_m)$  is derived by substituting (4.63) and (4.65) into (4.66).

(c) Using the definition of  $h_i(t_m, \theta_i, r_i)$ , and defining  $\overline{K_i^m} \triangleq \frac{(1+2\sigma^2 K_i^m)}{2\pi\sigma^2}$ , by substituting (4.67) in (4.68), the following is obtained

$$\begin{aligned} \hat{h}_i(t_m, \theta_i, r_i) &= \int_0^\infty \int_0^{2\pi} \left[ \cos(\omega_c(t_m - \tau_i(t_m)) + \theta_i) S(t_m - \tau_i(t_m)) r_i \overline{K_i^m} \exp\left(-r_i^2 \overline{K_i^m}\right) \right. \\ &\quad \times \exp\left(r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m))\right) \exp\left(-\frac{V_i^2(y^m)\sigma^2}{2 + 4\sigma^2 K_i^m}\right) \left. \right] d\theta_i dr_i. \end{aligned} \quad (\text{D.3})$$

Now, by writing  $\cos(\omega_c(t_m - \tau_i(t_m)) + \theta_i) = \cos(\omega_c(t_m - \tau_i(t_m)) + \gamma_i(y^m) + \theta_i - \gamma_i(y^m))$  and using the trigonometric identity  $\cos(A+B) = \cos(A)\cos(B) - \sin(A)\sin(B)$ , after some algebra, (D.3) transforms into

$$\begin{aligned} \hat{h}_i(t_m, \theta_i, r_i) &= S(t_m - \tau_i(t_m)) \cos(\omega_c(t_m - \tau_i(t_m)) + \gamma_i(y^m)) 2\pi \overline{K_i^m} \\ &\quad \times \exp\left(-\frac{V_i^2(y^m)\sigma^2}{2 + 4\sigma^2 K_i^m}\right) \int_0^\infty r_i^2 \exp\left(-2\pi r_i^2 \overline{K_i^m}\right) I_1\left(r_i V_i(y^m)\right) dr_i \end{aligned} \quad (\text{D.4})$$

where  $I_1(\cdot)$  is the modified first order Bessel function of the first order defined by  $I_1(x) \triangleq \frac{1}{2\pi} \int_{-\pi}^\pi \cos(\alpha) \exp(x \cos \alpha) d\alpha$ .

Next, by substituting  $\bar{\sigma}^{2,m}$  and  $\epsilon^m$  into (D.4), the following is obtained

$$\begin{aligned} \hat{h}_i(t_m, \theta_i, r_i) &= S(t_m - \tau_i(t_m)) \cos(\omega_c(t_m - \tau_i(t_m)) + \gamma_i(y^m)) \frac{1}{\bar{\sigma}^{2,m}} \exp\left(-\frac{\epsilon^{2,m}}{2\bar{\sigma}^{2,m}}\right) \\ &\quad \times \int_0^\infty r_i^2 \exp\left(-\frac{r_i^2}{2\bar{\sigma}^{2,m}}\right) I_1\left(\frac{r_i \epsilon^m}{\bar{\sigma}^{2,m}}\right) dr_i \\ &= S(t_m - \tau_i(t_m)) \cos(\omega_c(t_m - \tau_i(t_m)) + \gamma_i(y^m)) \epsilon^m \\ &\quad \times \int_0^\infty \frac{r_i^2}{\bar{\sigma}^{2,m} \epsilon^m} \exp\left(-\frac{(r_i^2 + \epsilon^{2,m})}{2\bar{\sigma}^{2,m}}\right) I_1\left(\frac{r_i \epsilon^m}{\bar{\sigma}^{2,m}}\right) dr_i. \end{aligned} \quad (\text{D.5})$$

Note that in (D.5) the fact that  $\frac{r_i^2}{\bar{\sigma}^{2,m}\epsilon^m} \exp\left(-\frac{(r_i^2+\epsilon^{2,m})}{2\bar{\sigma}^{2,m}}\right) I_1\left(\frac{r_i\epsilon^m}{\bar{\sigma}^{2,m}}\right)$  is the pdf of a noncentral chi-square distribution with 2 degrees of freedom and noncentrality parameter  $\epsilon^{2,m}$ , is being used, thus it integrates to 1. By taking this into consideration equation (4.69) is derived.

(d) The minimum least-square estimator  $\tilde{\theta}_i^*(t_m)$  is derived by substituting (4.67) into (4.70).

(e) Substituting (4.67) in (4.72), the following function is obtained

$$\begin{aligned}\tilde{r}_i^*(t_m) &= \int_0^\infty \int_0^{2\pi} r_i \frac{r_i(1+2\sigma^2 K_i^m)}{2\pi\sigma^2} \exp\left(-\frac{r_i^2(1+2\sigma^2 K_i^m)}{2\sigma^2}\right) \\ &\quad \times \exp\left(r_i V_i(y^m) \cos(\theta_i - \gamma_i(y^m))\right) \exp\left(-\frac{V_i^2(y^m)\sigma^2}{2+4\sigma^2 K_i^m}\right) d\theta_i dr_i \\ &= \exp\left(-\frac{V_i^2(y^m)\sigma^2}{2+4\sigma^2 K_i^m}\right) \int_0^\infty r_i \frac{r_i(1+2\sigma^2 K_i^m)}{\sigma^2} \\ &\quad \times \exp\left(-\frac{r_i^2(1+2\sigma^2 K_i^m)}{2\sigma^2}\right) I_0\left(r_i V_i(y^m)\right) dr_i.\end{aligned}\tag{D.6}$$

By substituting  $\bar{\sigma}^{2,m}$  and  $\epsilon^m$  into (D.6), then

$$\begin{aligned}\tilde{r}_i^*(t_m) &= \exp\left(-\frac{V_i^2(y^m)\sigma^2}{2+4\sigma^2 K_i^m}\right) \exp\left(\frac{\epsilon^{2,m}}{2\bar{\sigma}^{2,m}}\right) \\ &\quad \times \int_0^\infty r_i \frac{r_i}{\bar{\sigma}^{2,m}} \exp\left(-\frac{(r_i^2+\epsilon^{2,m})}{2\bar{\sigma}^{2,m}}\right) I_0\left(\frac{r_i\epsilon^m}{\bar{\sigma}^{2,m}}\right) dr_i \\ &= \exp\left(-\frac{V_i^2(y^m)\sigma^2}{2+4\sigma^2 K_i^m}\right) \exp\left(\frac{\epsilon^{2,m}}{2\bar{\sigma}^{2,m}}\right) E[r].\end{aligned}\tag{D.7}$$

Note that  $\frac{r_i}{\bar{\sigma}^{2,m}} \exp\left(-\frac{(r_i^2+\epsilon^{2,m})}{2\bar{\sigma}^{2,m}}\right) I_0\left(\frac{r_i\epsilon^m}{\bar{\sigma}^{2,m}}\right)$  is the pdf of a Ricean distribution random variable, hence the integral  $\int_0^\infty r_i \frac{r_i}{\bar{\sigma}^{2,m}} \exp\left(-\frac{r_i^2}{2\bar{\sigma}^{2,m}}\right) I_0\left(\frac{r_i\epsilon^m}{\bar{\sigma}^{2,m}}\right) dr_i$  is its first moment. Using the expression of the moments of a Rice distribution which are given by

$$E[r_i^k] = (2\bar{\sigma}^{2,m})^{k/2} \exp\left(-\frac{\epsilon^{2,m}}{2\bar{\sigma}^{2,m}}\right) \Gamma\left(1+\frac{k}{2}\right) {}_1F_1\left(\frac{2+k}{2}, 1; \frac{\epsilon^{2,m}}{2\bar{\sigma}^{2,m}}\right),$$

where  ${}_1F_1(\alpha, \beta; x)$  is the confluent hypergeometric function equation (4.73) is obtained.



## APPENDIX E

### DERIVATIONS OF EXAMPLES OF CHAPTER 5

This Appendix presents the derivations of the solutions of the Examples of Chapter 5.

#### Derivations of Example of Section 5.4.1 - Generalized MAP Estimator

First the Generalized MAP estimator  $\Phi_{MAP}$ , which is given by (5.27), is derived. Using (5.20), (5.21) and (5.22) the function  $\bar{L}(x, y^m)$  is given by

$$\begin{aligned}
 \bar{L}(x, y^m) &\triangleq \frac{\ell(x)}{s} + \log \frac{\mu(x, dy^m)}{dy^m} + \log \frac{dP_x(x)}{dx} \\
 &= \frac{x^T U x}{s} + \log \left( \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}} \right)^{m+1} \\
 &\quad - \sum_{i=0}^m (y_i - Hx)^T \frac{\Sigma_W^{-1}}{2} (y_i - Hx) \\
 &\quad + \log \left( \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_X|^{\frac{1}{2}}} \right) - x^T \frac{\Sigma_X^{-1}}{2} x \\
 &= \frac{x^T U x}{s} + \log \left( (2\pi)^{-\frac{d(m+1)}{2}} (2\pi)^{-\frac{n}{2}} |\Sigma_W|^{-\frac{1}{2}} |\Sigma_X|^{\frac{1}{2}} \right) \\
 &\quad - \sum_{i=0}^m y_i^T \frac{\Sigma_W^{-1}}{2} y_i - N x^T H^T \frac{\Sigma_W^{-1}}{2} H x \\
 &\quad + 2x^T H^T \frac{\Sigma_W^{-1}}{2} \sum_{i=0}^m y_i - x^T \frac{\Sigma_X^{-1}}{2} x.
 \end{aligned} \tag{E.1}$$

Next, the MAP estimator is computed by differentiating the above pay-off function and setting the derivative to zero as shown in (5.26). That is

$$\frac{\partial}{\partial x} \bar{L}(x, y^m) \Big|_{x=\Phi_{MAP}(y^m)} = 0 \quad (\text{E.2})$$

$$\begin{aligned} &\Rightarrow \frac{2(m+1)}{s} \Phi_{MAP}(y^m) - (m+1)H^T \Sigma_W^{-1} H \Phi_{MAP}(y^m) + H^T \Sigma_W^{-1} \sum_{i=0}^m y_i \\ &\quad - \Sigma_X^{-1} \Phi_{MAP}(y^m) = 0 \\ &\Rightarrow \left( \frac{2(m+1)}{s} - (m+1)H^T \Sigma_W^{-1} H - \Sigma_X^{-1} \right) \Phi_{MAP}(y^m) + H^T \Sigma_W^{-1} \sum_{i=0}^m y_i = 0 \\ &\Rightarrow \Phi_{MAP}(y^m) = \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{s} \right)^{-1} H^T \Sigma_W^{-1} \sum_{i=0}^m y_i. \end{aligned} \quad (\text{E.3})$$

Secondly, the mean square error, given by (5.30), is derived. The mean square error is given by (5.28), and can be written as

$$\begin{aligned} &E \left[ (X - \Phi_{MAP}(Y^m))^T (X - \Phi_{MAP}(Y^m)) \right] = \\ &tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} (x - \Phi_{MAP}(y^m))(x - \Phi_{MAP}(y^m))^T \mu(x, dy^m) dP_X(x) \quad (\text{E.4}) \\ &= tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} \left( xx^T + \Phi_{MAP}(y^m) \Phi_{MAP}^T(y^m) - 2x \Phi_{MAP}^T(y^m) \right) \mu(x, dy^m) dP_X(x) \\ &= tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} xx^T \mu(x, dy^m) dP_X(x) \\ &\quad - tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} 2x \Phi_{MAP}^T(y^m) \mu(x, dy^m) dP_X(x) \\ &\quad + tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} \Phi_{MAP}(y^m) \Phi_{MAP}^T(y^m) \mu(x, dy^m) dP_X(x). \end{aligned} \quad (\text{E.5})$$

Next, given (5.27)  $\Phi_{MAP}(y^m) = \Sigma_{\Phi_{MAP}} \sum_{i=0}^m y_i$ , where  $\Sigma_{\Phi_{MAP}} = \left( (m+1)H^T \Sigma_W^{-1} H + \Sigma_X^{-1} - \frac{2U}{s} \right)^{-1} H^T \Sigma_W^{-1}$ , (E.5) can also be written as

$$\begin{aligned} &tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} xx^T \mu(x, dy^m) dP_X(x) - tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} 2x \sum_{i=0}^m y_i^T \Sigma_{\Phi_{MAP}}^T \mu(x, dy^m) dP_X(x) \\ &+ tr \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} \Sigma_{\Phi_{MAP}} \sum_{i=0}^m y_i \sum_{i=0}^m y_i^T \Sigma_{\Phi_{MAP}}^T \mu(x, dy^m) dP_X(x) \\ &= tr(\Sigma_X) - tr \left( 2 \Sigma_{\Phi_{MAP}}^T \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} x \sum_{i=0}^m (Hx + w_i)^T \mu(x, dy^m) dP_X(x) \right) \end{aligned}$$

$$\begin{aligned}
 & +tr\left(\Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} \sum_{i=0}^m \sum_{j=0}^m y_i y_j^T \mu(x, dy^m) dP_X(x)\right) \\
 & = tr(\Sigma_X) - tr\left(2\Sigma_{\Phi_{MAP}}^T \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} (m+1)xx^T H^T \mu(x, dy^m) dP_X(x)\right) \\
 & +tr\left(\Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} \sum_{i=0}^m \sum_{j=0}^m (Hx + w_i)(Hx + w_j)^T \mu(x, dy^m) dP_X(x)\right) \\
 & = tr(\Sigma_X) - tr(2(m+1)H^T \Sigma_{\Phi_{MAP}}^T \Sigma_X) \\
 & +tr\left(\Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} \sum_{i=0}^m \sum_{j=0}^m (Hxx^T H^T + xHw_j^T + w_i x^T H^T + w_i w_j^T) \right. \\
 & \left. \times \mu(x, dy^m) dP_X(x)\right) \\
 & = tr(\Sigma_X - 2(m+1)H^T \Sigma_{\Phi_{MAP}}^T \Sigma_X) \\
 & +tr\left(\Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} ((m+1)^2 Hxx^T H^T) \mu(x, dy^m) dP_X(x)\right) \\
 & +tr\left(\Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \int_{\mathbb{R}^n} \int_{\mathbb{R}^{(m+1)d}} (m+1)w_i w_i^T \mu(x, dy^m) dP_X(x)\right) \\
 & = tr(\Sigma_X - 2(m+1)H^T \Sigma_{\Phi_{MAP}}^T \Sigma_X) + tr\left((m+1)^2 H^T \Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} H \Sigma_X \right. \\
 & \left. + (m+1)\Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \Sigma_W\right) \\
 & = tr\left(\Sigma_X - 2(m+1)H^T \Sigma_{\Phi_{MAP}}^T \Sigma_X + (m+1)^2 H^T \Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} H \Sigma_X \right. \\
 & \left. + (m+1)\Sigma_{\Phi_{MAP}}^T \Sigma_{\Phi_{MAP}} \Sigma_W\right).
 \end{aligned}$$

(E.6)

## Derivations of Example of Section 5.4.2 - Generalized ML Estimator

First the Generalized ML estimator  $\Phi_{ML}$ , which is given by (5.35), is derived. The MAP estimator is computed by differentiating the pay-off function,  $\bar{L}_2(x, y^m)$ , and setting the derivative to zero as shown in (5.34). By expanding (5.33),  $\bar{L}_2(x, y^m)$  is given by

$$\bar{L}_2(x, y^m) = \frac{x^T U x}{s} + \log \left( (2\pi)^{-\frac{d}{2}} |\Sigma_W|^{-\frac{1}{2}} \right)^{m+1} - \sum_{i=0}^m y_i^T \frac{\Sigma_W^{-1}}{2} y_i$$

$$-(m+1)x^T H^T \frac{\Sigma_W^{-1}}{2} Hx + x^T H^T \Sigma_W^{-1} \sum_{i=0}^m y_i. \quad (\text{E.7})$$

Next, the differentiation is performed.

$$\begin{aligned} & \frac{\partial}{\partial x} \bar{L}_2(x, y^m) \Big|_{x=\Phi_{ML}(y^m)} = 0 \\ \Rightarrow & \frac{2U}{s} \Phi_{ML}(y^m) - (m+1)H^T \Sigma_W^{-1} H \Phi_{ML}(y^m) + H^T \Sigma_W^{-1} \sum_{i=0}^m y_i = 0 \\ \Rightarrow & \left( \frac{2U}{s} - (m+1)H^T \Sigma_W^{-1} H \right) \Phi_{ML}(y^m) + H^T \Sigma_W^{-1} \sum_{i=0}^m y_i = 0 \\ \Rightarrow & \Phi_{ML}(y^m) = \left( (m+1)H^T \Sigma_W^{-1} H - \frac{2U}{s} \right)^{-1} H^T \Sigma_W^{-1} \sum_{i=0}^m y_i. \end{aligned} \quad (\text{E.8})$$

Secondly the mean square error, given by (5.37), is derived. The mean square error is given by (5.36), and can be written as (for  $X = x$  deterministic)

$$\begin{aligned} & E \left[ (x - \Phi_{ML}(Y^m))^T (x - \Phi_{ML}(Y^m)) \right] = \\ & tr \int_{\mathfrak{R}^{(m+1)d}} (x - \Phi_{ML}(y^m))(x - \Phi_{ML}(y^m))^T \mu(x, dy^m) \end{aligned} \quad (\text{E.9})$$

$$\begin{aligned} & = tr \int_{\mathfrak{R}^{(m+1)d}} \left( xx^T + \Phi_{ML}(y^m) \Phi_{ML}^T(y^m) - 2x \Phi_{ML}^T(y^m) \right) \mu(x, dy^m) \\ & = tr \int_{\mathfrak{R}^{(m+1)d}} xx^T \mu(x, dy^m) - tr \int_{\mathfrak{R}^{(m+1)d}} 2x \Phi_{ML}^T(y^m) \mu(x, dy^m) \\ & + tr \int_{\mathfrak{R}^{(m+1)d}} \Phi_{ML}(y^m) \Phi_{ML}^T(y^m) \mu(x, dy^m). \end{aligned} \quad (\text{E.10})$$

Given (5.35)  $\Phi_{ML}(y^m) = \Sigma_{\Phi_{ML}} \sum_{i=0}^m y_i$ , where  $\Sigma_{\Phi_{ML}} = \left( (m+1)H^T \Sigma_W^{-1} H - \frac{2U}{s} \right)^{-1} H^T \Sigma_W^{-1}$ , (E.10) can also be written as

$$\begin{aligned} & tr \int_{\mathfrak{R}^{(m+1)d}} xx^T \mu(x, dy^m) - tr \int_{\mathfrak{R}^{(m+1)d}} 2x \sum_{i=0}^m y_i^T \Sigma_{\Phi_{ML}}^T \mu(x, dy^m) \\ & + tr \int_{\mathfrak{R}^{(m+1)d}} \Sigma_{\Phi_{ML}} \sum_{i=0}^m y_i \sum_{i=0}^m y_i^T \Sigma_{\Phi_{ML}}^T \mu(x, dy^m) \\ & = tr(xx^T) - tr \left( 2 \Sigma_{\Phi_{ML}}^T \int_{\mathfrak{R}^{(m+1)d}} x \sum_{i=0}^m (Hx + w_i)^T \mu(x, dy^m) \right) \\ & + tr \left( \Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} \int_{\mathfrak{R}^{(m+1)d}} \sum_{i=0}^m \sum_{j=0}^m y_i y_j^T \mu(x, dy^m) \right) \end{aligned}$$



$$\begin{aligned}
 &= \text{tr}(xx^T) - \text{tr}\left(2\Sigma_{\Phi_{ML}}^T \int_{\mathfrak{R}^{(m+1)d}} (m+1)xx^T H^T \mu(x, dy^m)\right) \\
 &+ \text{tr}\left(\Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} \int_{\mathfrak{R}^{(m+1)d}} \sum_{i=0}^m \sum_{j=0}^m (Hx + w_i)(Hx + w_j)^T \mu(x, dy^m)\right) \\
 &= \text{tr}(xx^T) - \text{tr}(2(m+1)H^T \Sigma_{\Phi_{ML}}^T xx^T) \\
 &+ \text{tr}\left(\Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} \int_{\mathfrak{R}^{(m+1)d}} \sum_{i=0}^m \sum_{j=0}^m (Hxx^T H^T + xHw_j^T + w_i x^T H^T + w_i w_j^T) \mu(x, dy^m)\right) \\
 &= \text{tr}(xx^T - 2(m+1)H^T \Sigma_{\Phi_{ML}}^T xx^T) \\
 &+ \text{tr}\left(\Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} \int_{\mathfrak{R}^{(m+1)d}} ((m+1)^2 Hxx^T H^T) \mu(x, dy^m)\right) \\
 &+ \text{tr}\left(\Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} \int_{\mathfrak{R}^{(m+1)d}} (m+1)w_i w_i^T \mu(x, dy^m)\right) \\
 &= \text{tr}(xx^T - 2(m+1)H^T \Sigma_{\Phi_{ML}}^T xx^T) + \text{tr}\left((m+1)^2 H^T \Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} Hxx^T\right) \\
 &+ (m+1)\Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} \Sigma_W) \\
 &= \text{tr}\left(xx^T - 2(m+1)H^T \Sigma_{\Phi_{ML}}^T xx^T + (m+1)^2 H^T \Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} Hxx^T\right) \\
 &+ (m+1)\Sigma_{\Phi_{ML}}^T \Sigma_{\Phi_{ML}} \Sigma_W).
 \end{aligned} \tag{E.11}$$

Finally, the derivation of the *Cramér – Rao lower bound* is presented. First, equation (5.42) is derived.

$$\begin{aligned}
 \text{Var}[\Phi_{ML}(Y^m)] &= E\left[\left(\Phi_{ML}(Y^m) - E[\Phi_{ML}(Y^m)]\right)\left(\Phi_{ML}(Y^m) - E[\Phi_{ML}(Y^m)]\right)^T\right] \\
 &= E\left[\Sigma_{\Phi_{ML}}\left(\sum_{i=0}^m y_i - \sum_{i=0}^m E[y_i]\right)\left(\sum_{i=0}^m y_i - \sum_{i=0}^m E[y_i]\right)^T \Sigma_{\Phi_{ML}}^T\right] \\
 &= \Sigma_{\Phi_{ML}} \sum_{i=0}^m \text{Var}[y_i] \Sigma_{\Phi_{ML}}^T \\
 &= (m+1)\Sigma_{\Phi_{ML}} \Sigma_W \Sigma_{\Phi_{ML}}^T.
 \end{aligned} \tag{E.12}$$

Now, given  $E[\Phi_{ML}(y^m)] = \Sigma_{\Phi_{ML}}(m+1)Hx$ , then

$$V_X = \frac{\partial}{\partial x} E_X[\Phi_{ML}(y^m)] = \frac{\partial}{\partial x} \Sigma_{\Phi_{ML}}(m+1)Hx = \Sigma_{\Phi_{ML}}(m+1)H. \tag{E.13}$$

Furthermore,

$$\log \frac{\mu(x, dy^m)}{dy^m} = \log[(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}]^{-(m+1)} - \sum_{i=0}^m (y_i - Hx)^T \frac{\Sigma_W^{-1}}{2} (y_i - Hx)$$

$$\begin{aligned}
 &= \log[(2\pi)^{\frac{d}{2}} |\Sigma_W|^{\frac{1}{2}}]^{-(m+1)} - \sum_{i=0}^m y_i^2 \frac{\Sigma_W^{-1}}{2} y_i \\
 &\quad - (m+1) x^T H^T \frac{\Sigma_W^{-1}}{2} H x + x^T H^T \Sigma_W^{-1} \sum_{i=0}^m y_i.
 \end{aligned} \tag{E.14}$$

Next, the first derivative of (E.14) is given by

$$\frac{\partial}{\partial x} \log \frac{\mu(x, dy^m)}{dy^m} = -(m+1) H^T \Sigma_w^{-1} H x + H^T \Sigma_W^{-1} \sum_{i=0}^m y_i, \tag{E.15}$$

and the second derivative is given by

$$\frac{\partial^2}{\partial x^2} \log \frac{\mu(x, dy^m)}{dy^m} = -(m+1) H^T \Sigma_w^{-1} H. \tag{E.16}$$

Finally,

$$I_X = -E_X \left[ \frac{\partial^2}{\partial x^2} \log \frac{\mu(x, dy^m)}{dy} \right] = (m+1) H^T \Sigma_w^{-1} H \tag{E.17}$$

## LIST OF ACRONYMS

1. a.s. - almost surely
2. AVC - Arbitrarily Varying Channel
3. BME - Blind Minimax Estimator
4. EKF - Extended Kalman Filter
5. iid - independent and identically distributed
6. KL - Kullback-Leibler
7. LS - Least-Square
8. MAP - Maximum  $\hat{A}$  Posteriori
9. MIMO - Multiple-Input Multiple-Output
10. ML - Maximum Likelihood
11. MSE - Mean-Square-Error
12. MMAE - Minimum-Mean-Absolute-Error
13. MMSE - Minimum-Mean-Squared-Error
14. PDF - Probability Density Function
15. PSD - Power Spectral Density
16. RND -Radon Nikodym Derivative
17. RP - Random Process
18. RV - Random Variable
19. SNR - Signal-to-Noise Ratio



## NOTATION AND LIST OF SYMBOLS

1.  $N \triangleq \{1, 2, \dots\}$ .
2.  $N_0 \triangleq \{0, 1, 2, \dots\}$ .
3.  $N_0^m \triangleq \{0, 1, 2, \dots, m\}$ .
4.  $\mathbb{C}$ : set of complex numbers.
5.  $\mathbb{R}$ : set of real numbers.
6.  $\mathcal{Z}$ : set.
7.  $B^T$ : transpose of a matrix  $B$ .
8.  $|B|$ : determinant of a square matrix  $B$ .
9.  $\Omega$ : elementary outcomes of a random experiment.
10.  $\mathcal{F}$ :  $\sigma$ -algebra (or algebra) associated with a random experiment.
11.  $\mathcal{X}$ : Polish space.
12.  $\mathcal{Y}$ : Polish space.
13.  $\Sigma_{\mathcal{X}}$  and  $\Sigma_{\mathcal{Y}}$ :  $\sigma$ -algebras generated by  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively.
14.  $\mathcal{M}_1(\mathcal{X})$ : space of probability measures on  $\mathcal{X}$ .
15.  $\ell$ : sample pay-off.
16.  $E[\cdot]$ : expectation.
17.  $E[X|Y]$ : conditional expectation of  $X$  given  $Y$ .

18.  $BC(\mathcal{Y})$ : vector space of bounded continuous real valued functions defined on the Polish space  $\mathcal{Y}$
19.  $(BC(\mathcal{Y}))^*$ : topological dual of  $BC(\mathcal{Y})$ .
20.  $P_{X,Y}$ : nominal joint distribution of  $X$  and  $Y$ .
21.  $P_{X|Y}$ : conditional distribution of  $X$  given  $Y$ .
22.  $Q_{X,Y}$ : true joint distribution of  $X$  and  $Y$ .
23.  $Q_{X|Y}$ : true conditional distribution of  $X$  given  $Y$ .
24.  $P_X$ : marginal probability distribution of  $X$  associated with joint probability of  $X, Y$   $P_{X,Y}$ .
25.  $\mu, \eta, \nu$ : stochastic kernels.
26.  $\mathcal{P}$ : class of all stochastic kernel.
27.  $M_{rba}(\mathcal{Y})$ : Banach space of finitely additive regular bounded signed measures on  $\Sigma_{\mathcal{Y}}$ .
28.  $\Pi_{rba}(\mathcal{Y})$ : set of regular bounded finitely additive probability measures on  $\mathcal{Y}$ .
29.  $L_1(P_X, BC(\mathcal{Y}))$ : space of all  $P_X$  integrable functions defined on  $\mathcal{X}$  with values in  $BC(\mathcal{Y})$ .
30.  $H(P|Q)$ : relative entropy between the probability measures  $P$  with respect to  $Q$ .
31.  $\hat{X} = \Phi(Y)$ : estimate of  $X$  from the measurements of  $Y$ .
32.  $\Phi^*(Y)$ : best estimate of  $X$  from the measurements of  $Y$ .
33.  $\|\cdot\|_{\mathcal{X}}$ : norm associated with elements in  $\mathcal{X}$ .
34.  $\{\mathcal{G}_m^o\}$ :  $\sigma$ -field generated by the complete data  $\{x_0, x_1, \dots, x_m, y_0, y_1, \dots, y_m\}$ .
35.  $\{\mathcal{G}_m\}$ : complete filtration of  $\{\mathcal{G}_m^o\}$ .
36.  $\{\mathcal{Y}_m^o\}$ :  $\sigma$ -field generated by the incomplete data  $\{y_0, y_1, \dots, y_m\}$ .

37.  $\{\mathcal{Y}_m\}$ : complete filtration of  $\{\mathcal{Y}_m^o\}$ .
38.  $y^m$ : sequence  $\{y_0, \dots, y_m\}$  and similarly for other sequences.
39.  $\tilde{x}_m$ : estimate of the state  $x_m$  given  $\{\mathcal{Y}_m\}$ .
40.  $\alpha(\cdot)$ : unnormalized conditional probability distribution function.
41.  $\bar{\alpha}(\cdot)$ : unnormalized conditional probability density function.
42.  $\omega_c$ : carrier frequency.
43.  $\{\tau_i(t_k)\}$ : propagation delay.
44.  $\{r_i\}, \{\theta_i\}$ : attenuation and phase, respectively, of the signal received associated with  $i$ th path.
45.  $\pi_{r_0}(r_i)$ :  $\hat{a}$  priori probability density function of  $\{r_i\}$ , and similarly for other processes.
46.  $\hat{\Lambda}(t_m)$ : incomplete data likelihood ratio.





## BIBLIOGRAPHY

- [1] H. V. Poor, *An Introduction to Signal Detection and Estimation*. Springer-Verlag, 1994.
- [2] H. L. VanTrees, *Detection, Estimation and Modulation Theory-Part I*. Wiley, 1968.
- [3] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*. New York: McGraw-Hill, fourth ed., 2002.
- [4] T. Kailath and H. V. Poor, "Detection of stochastic processes," *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2230–2259, October 1998.
- [5] A. Jazwinski, *Stochastic Processes and Filtering Theory*. San Diego, CA: Academic, 1970.
- [6] S. Haykin, *Kalman Filtering and Neural Networks*. New York: Wiley, 2001.
- [7] B. Anderson and J. Moore, *Optimal Filtering*. Englewood Cliffs, NJ: Prentice-Hall, 1979.
- [8] A. Lapidoth and P. Narayan, "Reliable communication under channel uncertainty," *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2148–2177, October 1998.
- [9] V. A. H. Sayed, H. Nascimento, and F. A. M. Cipparrone, "A regularized robust design criterion for uncertain data," *SIAM Journal on Matrix Analysis and Applications*, vol. 23, no. 4, pp. 1120–1142, 2002.
- [10] A. H. Sayed, "A framework for state-space estimation with uncertain models," *IEEE Transactions on Automatic Control*, vol. 46, no. 7, pp. 998–1013, July 2001.

## BIBLIOGRAPHY

- [11] Y. Guo and B. C. Levy, "Robust MSE equalizer design for MIMO communication systems in the presence of model uncertainties," *IEEE Transactions on Signal Processing*, vol. 54, no. 5, pp. 1840–1852, May 2006.
- [12] A. G. Dabak and D. H. Johnson, "Geometrically based robust detection," in *Proc. Conf. Information Science and Systems (Baltimore, MD)*, pp. 73–77, March 1993.
- [13] R. H. Shumway and D. S. Stoffer, *Time Series Analysis and its Applications*. New York: Springer-Verlag, 2000.
- [14] S. Amari and H. Nagaoka, *Methods of Information Geometry*. American Mathematical Society, 2001.
- [15] S. A. Kassam and H. V. Poor, "Robust techniques for signal processing," *Proceedings of the IEEE*, vol. 73, no. 3, pp. 433–481, March 1985.
- [16] B. C. Levy and R. Nikoukhah, "Robust least-squares estimation with relative entropy constraint," *IEEE Transactions on Information Theory*, vol. 50, no. 1, pp. 89–104, January 2004.
- [17] M. C. Yovits and J. L. Jackson, "Linear filter optimization with game theory considerations," *IRE National Convention Record*, pt. 4, pp. 193–199, 1955.
- [18] A. G. Carlton and J. W. Follin, "Recent developments in fixed and adaptive filtering," in *Proc. 2nd AGARD Guided Missiles Seminar, AGARDograph21, NATO Advanced Group for R & D*, pp. 285–300, 1956.
- [19] P. J. Huber, "Robust estimation of a location parameter," *Math. Statist.*, vol. 35, pp. 73–101, 1964.
- [20] P. J. Huber, "A robust version of the probability ratio test," *Ann. Math. Statist.*, vol. 36, pp. 1753–1758, 1965.
- [21] W. L. Root, "Stability in signal detection problems," in *Proc. Symp. in Applied Mathematics*, vol. 16, pp. 247–263, 1964.
- [22] F. R. Hampel, "A general qualitative definition of robustness," *Ann. Math. Stat.*, vol. 42, pp. 1887–1896, 1971.

- [23] Z. Ben-Haim and Y. C. Eldar, "Blind minimax estimation," *IEEE Transactions on Information Theory*, vol. 53, no. 9, pp. 3145–3157, September 2007.
- [24] Y. C. Eldar, "Comparing between estimation approaches: Admissible and dominating linear estimators," *IEEE Transactions on Signal Processing*, vol. 54, no. 5, pp. 1689–1702, May 2006.
- [25] Y. C. Eldar, A. Ben-Tal, and A. Nemirovski, "Robust mean-squared error estimation in the presence of model uncertainties," *IEEE Transactions on Signal Processing*, vol. 53, no. 1, pp. 168–181, January 2005.
- [26] Z. Ben-Haim and Y. C. Eldar, "Minimax estimators dominating the least-squares estimator," in *Proc. Int. Conf. Acoustics, Speech and Signal Processing (ICASSP 2005)*, vol. IV, pp. 53–56, March 2005.
- [27] J. H. Marton, V. Krishnamurthy, and H. V. Poor, "James-Stein state filtering algorithms," *IEEE Transactions on Signal Processing*, vol. 46, no. 9, pp. 2431–2447, September 1998.
- [28] Y. C. Eldar and N. Merhav, "A competitive minimax approach to robust estimation of random parameters," *IEEE Transactions on Signal Processing*, vol. 52, no. 7, pp. 1931–1946, July 2004.
- [29] R. K. Boel, M. R. James, and I. R. Petersen, "Robustness and risk-sensitive filtering," *IEEE Transactions on Automatic Control*, vol. 47, no. 3, pp. 451–461, March 2002.
- [30] K. Vastola and H. V. Poor, "Robust Wiener-Kolmogorov theory," *IEEE Transactions on Information Theory*, vol. 30, no. 2, pp. 316–326, 1984.
- [31] D. Guo, S. Shamai, and S. Verdú, "Mutual information and minimum mean-square error in gaussian channels," *IEEE Transactions on Information Theory*, vol. 51, no. 4, pp. 1261–1283, April 2005.
- [32] Y. Wu and S. Verdú, "MMSE dimension," in *Proc. of 2010 IEEE International Symposium in Information Theory (ISIT 2010)*, pp. 1463–1467, June 2010.

## BIBLIOGRAPHY

- [33] Y. Wu and S. Verdú, "Functional properties of MMSE," in *Proc. of 2010 IEEE International Symposium in Information Theory (ISIT 2010)*, pp. 1453–1457, June 2010.
- [34] T. Feng, T. R. Field, and S. Haykin, "Stochastic differential equation theory applied to wireless channels," *IEEE Transactions on Communications*, vol. 55, no. 8, pp. 1478–1483, August 2007.
- [35] K. E. Baddour and N. C. Beaulieu, "Autoregressive modeling for fading channel simulation," *IEEE Transactions on Wireless Communications*, vol. 4, no. 4, pp. 1650–1662, July 2005.
- [36] R. Elliott, L. Aggoun, and J. Moore, *Hidden Markov Models: Estimation and Control*. Springer-Verlag, Berlin-Heidelberg-New York 29, 1994.
- [37] N. U. Ahmed, *Dynamic Systems and Control with Applications*. World Scientific, Singapore, 2006.
- [38] F. Rezaei, C. D. Charalambous, and N. Ahmed, "Optimization of stochastic uncertain systems with variational norm constraints," in *46th IEEE Conference on Decision and Control*, pp. 2159–2163, December 2007.
- [39] K. Fan, "Minmax theorems," in *Proc. Nat. Acad. Sci.*, pp. 42–47, 1953.
- [40] D. G. Luenberger, *Optimization by Vector Space Method*. John Wiley, 1969.
- [41] C. D. Charalambous, "Lecture notes on random processes." Dept. of ECE, University of Cyprus., 2008.
- [42] E. Telatar, "Capacity of multi-antenna gaussian channels," *European Trans. Telecomm.*, vol. 10, no. 6, pp. 585–595, 1999.
- [43] C. D. Charalambous and F. Rezaei, "Stochastic uncertain systems subject to relative entropy constraints: Induced norms and monotonicity properties of min-max games," *IEEE Transactions on Automatic Control*, vol. 52, no. 4, pp. 647–663, May 2007.

- [44] P. D. Pra, L. Meneghini, and W. Runggaldier, "Some connections between stochastic control and dynamic games," *Mathematics of Control Signals, and Systems*, vol. 9, pp. 303–326, 1996.
- [45] N. Dunford and J. T. Schwartz, *Linear Operators, Part I: General Theory*. Interscience Publishers, Inc., New York, 1958.
- [46] A. I. Tulcea and C. I. Tulcea, *Topics in the Theory of Lifting*. Springer Verlag, Berlin, Heidelberg, New York, 1969.
- [47] P. Dupuis and R. S. Ellis, *A Weak Convergence Approach to the Theory of Large Deviations*. John Wiley, 1997.
- [48] G. J. Foschini, "Layered space-time architecture for wireless communications in a fading environment," *Bell Labs Tech. J.*, vol. 1, no. 2, pp. 41–59, 1996.
- [49] G. J. Foschini and M. J. Gans, "On the limit of wireless communications in a fading environment when using multiple antennas," *Wireless Personal Commun.*, vol. 6, no. 3, pp. 311–335, 1998.
- [50] C. D. Charalambous and Y. Socratous, "Nonlinear estimation for a class of systems," in *Proc. of the 2006 IEEE Symposium in Information Theory*, pp. 841–845, July 2006.
- [51] D. E. Asraf and M. G. Gustafsson, "An analytical series expansion to the problem of noncoherent detection," *IEEE Transactions on Information Theory*, vol. 50, no. 12, pp. 3369–3375, December 2004.
- [52] J. G. Proakis, *Digital Communications*. New York: McGraw-Hill, 4th ed., 2001.
- [53] C. D. Charalambous, A. Nejad, and D. Makrakis, "Coherent and noncoherent wireless channel estimation via ML/MAP and EM algorithm," in *Proc. of the 21st Biennial Symposium on Communications*, pp. 241–244, June 2002.
- [54] T. T. Georgiou and A. Lindquist, "Kullback-Leibler approximations of spectral density functions," *IEEE Transactions on Information Theory*, vol. 49, no. 11, pp. 2910–2917, November 2003.

## BIBLIOGRAPHY

- [55] M. S. Pinsker, *Information and Information Stability of Random Variables and Processes*. Holden-Day, San Francisco, 1964.
- [56] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge, U.K.: Cambridge Uni. Press, 1985.

Yiannis Socratous

# CURRICULUM VITAE

## **Yiannis Socratous**

E-mail: yianniss@ucy.ac.cy, ysocratous@gmail.com

### **Education**

- 2004 - 2010 PhD in Electrical Engineering,  
Electrical and Computer Engineering Department,  
University of Cyprus.
- 2001 - 2002 Master of Science in Mobile and Personal Communications,  
King's College London, UK. Grade: with Distinction.
- 1998 - 2001 Bachelor of Engineering in Electrical and Electronic Engineering,  
University of Plymouth, UK. Grade: First Class Honours.

### **Publications**

#### **Journals**

1. Y. Socratous, F. Rezaei, and C. D. Charalambous, "Nonlinear estimation for a class of systems", *IEEE Transactions on Information Theory*, vol. 55, no. 4, pp. 1930-1938, April 2009.

## Conferences

1. C. D. Charalambous and Y. Socratous, "Nonlinear estimation for a class of systems", in *Proc. of IEEE International Symposium on Information Theory*, pp. 841-845, 2006.
2. Y. Socratous, C. D. Charalambous and C. N. Georghiades, "Least-square estimation for nonlinear systems with applications to phase and envelope estimation in wireless fading channels", in *Proc. of the 47th IEEE Conference on Decision and Control*, pp. 4928-4932, 2008.
3. Y. Socratous, C. D. Charalambous and C. N. Georghiades, "Robust estimation with applications to phase and envelope estimation in frequency selective wireless fading channels", in *Proc. of IEEE International Symposium on Information Theory*, pp. 1238-1242, 2008.

## Pending Publications

### Journals

1. "Non-Coherent Detection and Estimation".
2. "Robust Least-square estimation for a class of systems".