
BAYESIAN ESTIMATION OF TIME-VARYING SYSTEMS:

Discrete-Time Systems

Simo Särkkä

Contents

Contents	i
1 Introduction	1
1.1 Why Bayesian Approach?	1
1.2 What is Optimal Filtering?	2
1.2.1 Applications of Optimal Filtering	2
1.2.2 Origins of Bayesian Optimal Filtering	6
1.2.3 Optimal Filtering and Smoothing as Bayesian Inference	7
1.2.4 Algorithms for Optimal Filtering and Smoothing	9
2 From Bayesian Inference to Bayesian Optimal Filtering	13
2.1 Bayesian Inference	13
2.1.1 Philosophy of Bayesian Inference	13
2.1.2 Connection to Maximum Likelihood Estimation	13
2.1.3 The Building Blocks of Bayesian Models	14
2.1.4 Bayesian Point Estimates	16
2.1.5 Numerical Methods	17
2.2 Batch and Recursive Estimation	19
2.2.1 Batch Linear Regression	19
2.2.2 Recursive Linear Regression	21
2.2.3 Batch vs. Recursive Estimation	22
2.3 Towards Bayesian Filtering	24
2.3.1 Drift Model for Linear Regression	24
2.3.2 Kalman Filter	26
3 Optimal Filtering	31
3.1 Formal Filtering Equations and Exact Solutions	31
3.1.1 Discrete-Time Probabilistic State Space Models	31
3.1.2 Optimal Filtering Equations	33
3.1.3 Kalman Filter	35
3.2 Gaussian Approximation Based Filtering	39
3.2.1 Linearization of Non-Linear Transforms	39
3.2.2 Extended Kalman Filter	42

3.2.3	Statistical Linearization of Non-Linear Transforms	46
3.2.4	Statistically Linearized Filter	48
3.2.5	Unscented Transform	49
3.2.6	Unscented Kalman Filter	54
3.2.7	Gaussian Moment Matching	56
3.2.8	Gaussian Assumed Density Filter	57
3.3	Monte Carlo Approximations	58
3.3.1	Principles and Motivation of Monte Carlo	58
3.3.2	Importance Sampling	59
3.4	Particle Filtering	60
3.4.1	Sequential Importance Sampling	60
3.4.2	Sequential Importance Resampling	61
3.4.3	Rao-Blackwellized Particle Filter	64
4	Optimal Smoothing	69
4.1	Formal Equations and Exact Solutions	69
4.1.1	Optimal Smoothing Equations	69
4.1.2	Discrete-Time Rauch-Tung-Striebel Smoother	70
4.2	Gaussian Approximation Based Smoothing	73
4.2.1	Discrete-Time Extended Rauch-Tung-Striebel Smoother	73
4.2.2	Statistically Linearized RTS Smoother	75
4.2.3	Unscented Rauch-Tung-Striebel Smoother	76
4.2.4	Gaussian Assumed Density RTS Smoother	78
4.3	Monte Carlo Based Smoothers	79
4.3.1	Sequential Importance Resampling Smoother	79
4.3.2	Rao-Blackwellized Particle Smoother	80
A	Additional Material	81
A.1	Properties of Gaussian Distribution	81
	References	83

Chapter 1

Introduction

1.1 Why Bayesian Approach?

The mathematical treatment of the models and algorithms in this document is Bayesian, which means that all the results are treated as being approximations to certain probability distributions or their parameters. Probability distributions are used for modeling both the uncertainties in the models and for modeling the physical randomness. The theory of non-linear optimal filtering is formulated in terms of Bayesian inference and both the classical and recent filtering algorithms are derived using the same Bayesian notation and formalism.

The reason for selecting the Bayesian approach is more a practical engineering than a philosophical decision. It simply is easier to develop a consistent, practically applicable theory of recursive inference under Bayesian philosophy than under, for example, least squares or maximum likelihood philosophy. Another useful consequence of selecting the Bayesian approach is that least squares, maximum likelihood and many other philosophically different results can be obtained as special cases or re-interpretations of the Bayesian results. Of course, quite often the same thing applies also other way around.

Modeling uncertainty as randomness is a very “engineering” way of modeling the world. It is exactly the approach also chosen in statistical physics as well as in financial analysis. Also the Bayesian approach to optimal filtering is far from new (see, e.g., Ho and Lee, 1964; Lee, 1964; Jazwinski, 1966; Stratonovich, 1968; Jazwinski, 1970), because the theory already existed at the same time the seminal article of Kalman (1960b) was published. The Kalman filter was derived from the least squares point of view, but the non-linear filtering theory has been Bayesian from the beginning (see, e.g., Jazwinski, 1970).

One should not take the Bayesian way of modeling unknown parameters as random variables too literally. It does not imply that one believes that there really is something random in the parameters - it is just a convenient way of representing uncertainty under the same formalism that is used for representing randomness. Also random or stochastic processes appearing in the mathematical models are

not necessarily really random in physical sense, but instead, the randomness is just a mathematical trick for taking into account the uncertainty in a dynamic phenomenon.

But it does not matter if the randomness is interpreted as physical randomness or as a representation of uncertainty, as long as the randomness based models succeed in modeling the real world. In the above engineering philosophy the controversy between so called “frequentists” and “Bayesians” is simply silly – it is quite much equivalent to the unnecessary controversy about interpretations of quantum mechanics, that is, whether, for example, the Copenhagen interpretation or several worlds implementation is the correct one. The philosophical interpretation does not matter as long as we get meaningful predictions from the theory.

1.2 What is Optimal Filtering?

Optimal filtering refers to the methodology that can be used for estimating the state of a time-varying system, which is indirectly observed through noisy measurements. The *state* of the system refers to the collection of dynamic variables such as position, velocities and accelerations or orientation and rotational motion parameters, which describe the physical state of the system. The *noise* in the measurements refers to a noise in the sense that the measurements are uncertain, that is, even if we knew the true system state the measurements would not be deterministic functions of the state, but would have certain distribution of possible values. The time evolution of the state is modeled as a dynamic system, which is perturbed by a certain *process noise*. This noise is used for modeling the uncertainties in the system dynamics and in most cases the system is not truly stochastic, but the stochasticity is only used for representing the model uncertainties.

1.2.1 Applications of Optimal Filtering

Phenomena, which can be modeled as time varying systems of the above type are very common in engineering applications. These kind of models can be found, for example, in navigation, aerospace engineering, space engineering, remote surveillance, telecommunications, physics, audio signal processing, control engineering, finance and several other fields. Examples of such applications are the following:

- *Global positioning system (GPS)* (Kaplan, 1996) is a widely used satellite navigation system, where the GPS receiver unit measures arrival times of signals from several GPS satellites and computes its position based on these measurements. The GPS receiver typically uses an extended Kalman filter or some other optimal filtering algorithm for computing the position and velocity such that the measurements and the assumed dynamics (laws of physics) are taken into account. Also the ephemeris information, which is the satellite reference information transmitted from the satellites to the GPS receivers is typically generated using optimal filters.

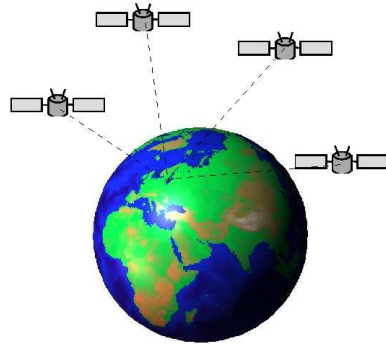


Figure 1.1: In GPS system, the measurements are time delays of satellite signals and the optimal filter (e.g., EKF) computes the position and the accurate time.

- *Target tracking* (Bar-Shalom et al., 2001) refers to the methodology, where a set of sensors such as active or passive radars, radio frequency sensors, acoustic arrays, infrared sensors and other types of sensors are used for determining the position and velocity of a remote target. When this tracking is done continuously, the dynamics of the target and measurements from the different sensors are most naturally combined using an optimal filter. The target in this (single) target tracking case can be, for example, a robot, a satellite, a car or an airplane.

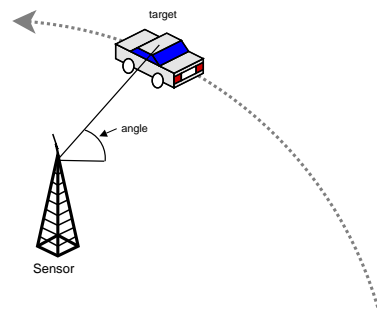


Figure 1.2: In target tracking, a sensor generates measurements (e.g., angle measurements) of the target, and the purpose is to determine the target trajectory.

- *Multiple target tracking* (Bar-Shalom and Li, 1995; Blackman and Popoli, 1999; Stone et al., 1999) systems are used for remote surveillance in the cases, where there are multiple targets moving at the same time in the same geographical area. This arises the concept of data association (which measurement was from which target?) and the problem of estimating the number

of targets. Multiple target tracking systems are typically used in remote surveillance for military purposes, but possible civil applications are, for example, monitoring of car tunnels, automatic alarm systems and people tracking in buildings.

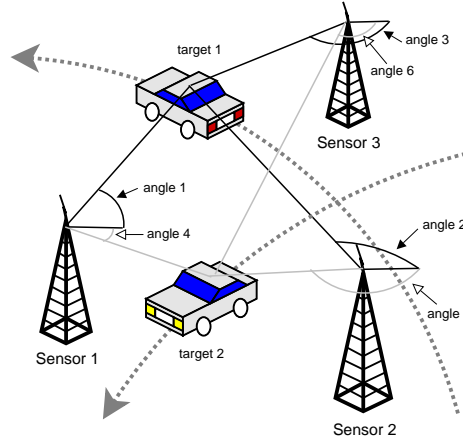


Figure 1.3: In multiple target tracking the data association problem has to be solved, which means that it is impossible to know without any additional information, which target produced which measurement.

- *Inertial navigation* (Titterton and Weston, 1997; Grewal et al., 2001) uses inertial sensors such as accelerometers and gyroscopes for computing the position and velocity of a device such as a car, an airplane or a missile. When the inaccuracies in sensor measurements are taken into account the natural way of computing the estimates is by using an optimal filter. Also in sensor calibration, which is typically done in time varying environment optimal filters are often applied.
- *Integrated inertial navigation* (Grewal et al., 2001; Bar-Shalom et al., 2001) combines the good sides of unbiased but inaccurate sensors, such as altimeters and landmark trackers, and biased but locally accurate inertial sensors. Combining of these different sources of information is most naturally performed using an optimal filter such as the extended Kalman filter. This kind of approach was used, for example, in the guidance system of Apollo 11 lunar module (Eagle), which landed on the moon in 1969.
- *GPS/INS navigation* (Grewal et al., 2001; Bar-Shalom et al., 2001) is a form of integrated inertial navigation, where the inertial sensors are combined with a GPS receiver unit. In GPS/INS navigation system the short term fluctuations of the GPS can be compensated with the inertial sensors and the inertial sensor biases can be compensated with the GPS receiver. An

additional advantage of this approach is that it is possible to temporarily switch to pure inertial navigation, when the GPS receiver is unable to compute its position (i.e., has no fix) for some reason. This happens, for example, indoors, in tunnels and in other cases when there is no direct line-of-sight between the GPS receiver and the satellites.

- *Spread of infectious diseases* (Anderson and May, 1991) can often be modeled as differential equations for the number of susceptible, infected and recovered/dead individuals. When uncertainties are induced into the dynamic equations, and when the measurements are not perfect, the estimation of the spread of the disease can be formulated as an optimal filtering problem.
- *Biological processes* (Murray, 1993) such as population growth, predator-prey models and several other dynamic processes in biology can also be modeled as (stochastic) differential equations. The estimation of the states of these processes from inaccurate measurements can be formulated as an optimal filtering problem.
- *Telecommunications* is also a field where optimal filters are traditionally used. For example, optimal receivers, signal detectors and phase locked loops can be interpreted to contain optimal filters (Van Trees, 1968, 1971) as components. Also the celebrated Viterbi algorithm (Viterbi, 1967) can be interpreted as a combination of optimal filtering and optimal smoothing of the underlying hidden Markov model.
- *Audio signal processing* applications such as audio restoration (Godsill and Rayner, 1998) and audio signal enhancement (Fong et al., 2002) often use TVAR (time varying autoregressive) models as the underlying audio signal models. These kind of models can be efficiently estimated using optimal filters and smoothers.
- *Stochastic optimal control* (Maybeck, 1982b; Stengel, 1994) considers control of time varying stochastic systems. Stochastic controllers can typically be found in, for example, airplanes, cars and rockets. The optimality, in addition to the statistical optimality, means that control signal is constructed to minimize a performance cost, such as expected time to reach a predefined state, the amount of fuel consumed or average distance from a desired position trajectory. Optimal filters are typically used for estimating the states of the stochastic system and a deterministic optimal controller is constructed independently from the filter such that it uses the estimate of the filter as the known state. In theory, the optimal controller and optimal filter are not completely decoupled and the problem of constructing optimal stochastic controllers is far more challenging than constructing optimal filters and (deterministic) optimal controllers separately.

- *Learning systems* or adaptive systems can often be mathematically formulated in terms of optimal filters. The theory of stochastic differential equations has close relationship with Bayesian non-parametric modeling, machine learning and neural network modeling (MacKay, 1998; Bishop, 1995). Methods, which are similar to the data association methods in multiple target tracking are also applicable to on-line adaptive classification (Andrieu et al., 2002).
- *Physical systems* which are time varying and measured through unideal sensors can sometimes be formulated as stochastic state space models, and the time evolution of the system can be estimated using optimal filters (Kaipio and Somersalo, 2005). In Vauhkonen (1997) and more recently, for example, in Pikkarainen (2005) optimal filtering is applied to Electrical Impedance Tomography (EIT) problem in time varying setting.

1.2.2 Origins of Bayesian Optimal Filtering

The roots of Bayesian analysis of time dependent behavior are in the optimal linear filtering. The idea of constructing mathematically optimal recursive estimators was first presented for linear systems due to their mathematical simplicity and the most natural optimality criterion in both mathematical and modeling point of view was the least squares optimality. For linear systems the optimal Bayesian solution (with MMSE utility) coincides with the least squares solution, that is, the optimal least squares solution is exactly the posterior mean.

The history of optimal filtering starts from the *Wiener filter* (Wiener, 1950), which is a spectral domain solution to the problem of least squares optimal filtering of stationary Gaussian signals. The Wiener filter is still important in communication applications (Proakis, 2001), digital signal processing (Hayes, 1996) and image processing (Rafael C. Gonzalez, 2008). The disadvantages of the Wiener filter are that it can only be applied to stationary signals and that the construction of a Wiener filter is often mathematically demanding and these mathematics cannot be avoided (i.e., made transparent). Due to the demanding mathematics the Wiener filter can only be applied to simple low dimensional filtering problems.

The success of optimal linear filtering in engineering applications is mostly due to the seminal article of Kalman (1960b), which describes the recursive solution to the optimal discrete-time (sampled) linear filtering problem. The reason to the success is that the *Kalman filter* can be understood and applied with very much lighter mathematical machinery than the Wiener filter. Also, despite its mathematical simplicity, the Kalman filter (or actually the Kalman-Bucy filter; Kalman and Bucy, 1961) contains the Wiener filter as its limiting special case.

In the early stages of its history, the Kalman filter was soon discovered to belong to the class of Bayesian estimators (Ho and Lee, 1964; Lee, 1964; Jazwinski, 1966, 1970). An interesting historical detail is that while Kalman and Bucy were formulating the linear theory in the United States, Stratonovich was doing the

pioneering work on the probabilistic (Bayesian) approach in Russia (Stratonovich, 1968; Jazwinski, 1970).

As discussed in the book of West and Harrison (1997), in the sixties, Kalman filter like recursive estimators were also used in the Bayesian community and it is not clear whether the theory of Kalman filtering or the theory of *dynamic linear models* (DLM) was the first. Although these theories were originally derived from slightly different starting points, they are equivalent. Because of Kalman filter's useful connection to the theory and history of stochastic optimal control, this document approaches the Bayesian filtering problem from the Kalman filtering point of view.

Although the original derivation of the *Kalman filter* was based on the least squares approach, the same equations can be derived from the pure probabilistic Bayesian analysis. The Bayesian analysis of Kalman filtering is well covered in the classical book of Jazwinski (1970) and more recently in the book of Bar-Shalom et al. (2001). Kalman filtering, mostly because of its least squares interpretation, has widely been used in stochastic optimal control. A practical reason to this is that the inventor of the Kalman filter, Rudolph E. Kalman, has also made several contributions (Kalman, 1960a) to the theory of *linear quadratic Gaussian* (LQG) regulators, which are fundamental tools of stochastic optimal control (Stengel, 1994; Maybeck, 1982b).

1.2.3 Optimal Filtering and Smoothing as Bayesian Inference

Optimal Bayesian filtering (see, e.g. Jazwinski, 1970; Bar-Shalom et al., 2001; Doucet et al., 2001; Ristic et al., 2004) considers statistical inversion problems, where the unknown quantity is a vector valued time series $(\mathbf{x}_1, \mathbf{x}_2, \dots)$ which is observed through noisy measurements $(\mathbf{y}_1, \mathbf{y}_2, \dots)$ as illustrated in the Figure 1.4. An example of this kind of time series is shown in the Figure 1.5. The process shown is actually a discrete-time noisy resonator with a known angular velocity. The state $\mathbf{x}_k = (x_k \dot{x}_k)^T$ is two dimensional and consists of the position of the resonator x_k and its time derivative \dot{x}_k . The measurements y_k are scalar observations of the resonator position (signal) and they are corrupted by measurement noise.

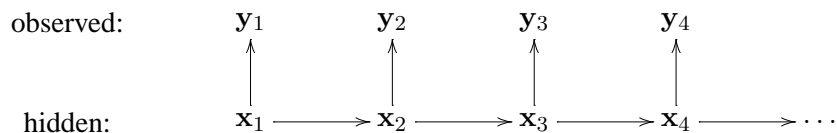


Figure 1.4: In discrete-time filtering a sequence of hidden states \mathbf{x}_k is indirectly observed through noisy measurements \mathbf{y}_k .

The purpose of the *statistical inversion* at hand is to estimate the hidden states $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ given the observed measurements $\{\mathbf{y}_1, \dots, \mathbf{y}_T\}$, which means that

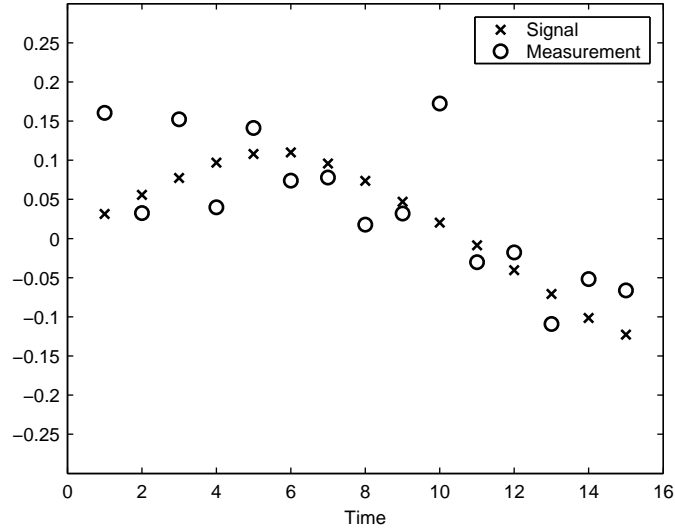


Figure 1.5: An example of time series, which models a discrete-time resonator. The actual resonator state (signal) is hidden and only observed through the noisy measurements.

in the Bayesian sense (Bernardo and Smith, 1994; Gelman et al., 1995) all we have to do is to compute the joint posterior distribution of all the states given all the measurements. This can be done by straightforward application of the Bayes' rule:

$$p(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{y}_1, \dots, \mathbf{y}_T) = \frac{p(\mathbf{y}_1, \dots, \mathbf{y}_T | \mathbf{x}_1, \dots, \mathbf{x}_T) p(\mathbf{x}_1, \dots, \mathbf{x}_T)}{p(\mathbf{y}_1, \dots, \mathbf{y}_T)}, \quad (1.1)$$

where

- $p(\mathbf{x}_1, \dots, \mathbf{x}_T)$, is the prior defined by the dynamic model,
- $p(\mathbf{y}_1, \dots, \mathbf{y}_T | \mathbf{x}_1, \dots, \mathbf{x}_T)$ is the likelihood model for the measurements,
- $p(\mathbf{y}_1, \dots, \mathbf{y}_T)$ is the normalization constant defined as

$$p(\mathbf{y}_1, \dots, \mathbf{y}_T) = \int p(\mathbf{y}_1, \dots, \mathbf{y}_T | \mathbf{x}_1, \dots, \mathbf{x}_T) p(\mathbf{x}_1, \dots, \mathbf{x}_T) d(\mathbf{x}_1, \dots, \mathbf{x}_T). \quad (1.2)$$

Unfortunately, this full posterior formulation has the serious disadvantage that each time we obtain a new measurement, the full posterior distribution would have to be recomputed. This is particularly a problem in dynamic estimation (which is exactly the problem we are solving here!), because there measurements are typically obtained one at a time and we would want to compute the best possible estimate after each measurement. When number of time steps increases, the dimensionality of the full posterior distribution also increases, which means that the computational

complexity of a single time step increases. Thus after a sufficient number of time steps the computations will become intractable, independently of available computational resources. Without additional information or harsh approximations, there is no way of getting over this problem in the full posterior computation.

However, the above problem only arises when we want to compute the *full* posterior distribution of the states at each time step. If we are willing to relax this a bit and be satisfied with selected marginal distributions of the states, the computations become order of magnitude lighter. In order to achieve this, we also need to restrict the class of dynamic models into probabilistic Markov sequences, which as a restriction sounds more restrictive than it really is. The model for the states and measurements will be assumed to be of the following type:

- **Initial distribution** specifies the *prior distribution* $p(\mathbf{x}_0)$ of the hidden state \mathbf{x}_0 at initial time step $k = 0$.
- **Dynamic model** models the system dynamics and its uncertainties as a *Markov sequence*, defined in terms of the transition distribution $p(\mathbf{x}_k | \mathbf{x}_{k-1})$.
- **Measurement model** models how the measurement \mathbf{y}_k depends on the current state \mathbf{x}_k . This dependence is modeled by specifying the distribution of the measurement given the state $p(\mathbf{y}_k | \mathbf{x}_k)$.

Because computing the full joint distribution of the states at all time steps is computationally very inefficient and unnecessary in real-time applications, in *optimal (Bayesian) filtering* the following marginal distributions are considered instead:

- *Filtering distributions* are the marginal distributions of *the current state* \mathbf{x}_k given *the previous measurements* $\{\mathbf{y}_1, \dots, \mathbf{y}_k\}$:

$$p(\mathbf{x}_k | \mathbf{y}_1, \dots, \mathbf{y}_k), \quad k = 1, \dots, T. \quad (1.3)$$

- *Prediction distributions* are the marginal distributions of the future states, n steps after the current time step:

$$p(\mathbf{x}_{k+n} | \mathbf{y}_1, \dots, \mathbf{y}_k), \quad k = 1, \dots, T, \quad n = 1, 2, \dots, \quad (1.4)$$

- *Smoothing distributions* are the marginal distributions of the states \mathbf{x}_k given a certain interval $\{\mathbf{y}_1, \dots, \mathbf{y}_T\}$ of measurements with $T > k$:

$$p(\mathbf{x}_k | \mathbf{y}_1, \dots, \mathbf{y}_T), \quad k = 1, \dots, T. \quad (1.5)$$

1.2.4 Algorithms for Optimal Filtering and Smoothing

There exists a few classes of filtering and smoothing problems which have closed form solutions:

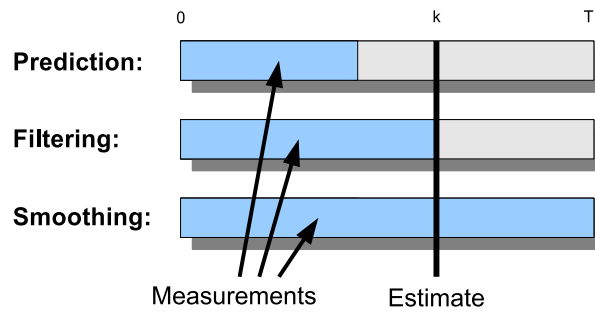


Figure 1.6: State estimation problems can be divided into optimal prediction, filtering and smoothing depending on the time span of measurements available with respect to the estimated state time span.

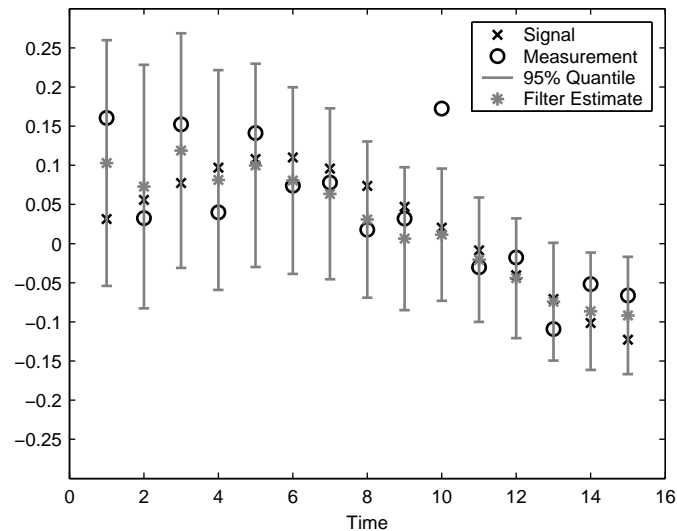


Figure 1.7: The result of computing the filtering distributions for the discrete-time resonator model. The *estimates* are the posterior means of the filtering distributions and the quantiles are the 95% quantiles of the filtering distributions.

- *Kalman filter* (KF) is a closed form solution to the discrete linear filtering problem. Due to linear Gaussian model assumptions the posterior distribution is exactly Gaussian and no numerical approximations are needed.
- *Rauch-Tung-Striebel smoother* (RTSS) is the corresponding closed form smoother to linear Gaussian state space models.

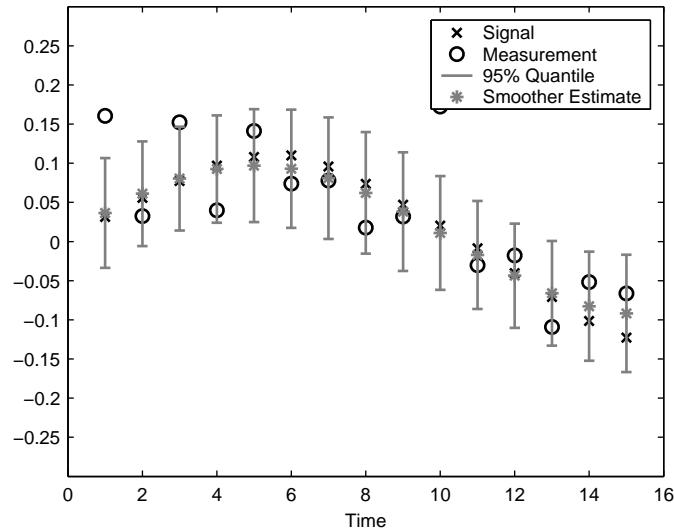


Figure 1.8: The result of computing the smoothing distributions for the discrete-time resonator model. The *estimates* are the posterior means of the smoothing distributions and the quantiles are the 95% quantiles of the smoothing distributions. The smoothing distributions are actually the marginal distributions of the full state posterior distribution.

- *Grid filters and smoothers*, are solutions to Markov models with finite state spaces.

But because the Bayesian optimal filtering and smoothing equations are generally computationally intractable, many kinds of numerical approximation methods have been developed:

- *Extended Kalman filter (EKF)* approximates the non-linear and non-Gaussian measurement and dynamic models by linearization, that is, by forming a Taylor series expansion on the nominal (or Maximum a Posteriori, MAP) solution. This results in Gaussian approximation to the filtering distribution.
- *Extended Rauch-Tung-Striebel smoother (ERTSS)* is the approximate non-linear smoothing algorithm corresponding to EKF.
- *Unscented Kalman filter (UKF)* approximates the propagation of densities through the non-linearities of measurement and noise processes by *unscented transform*. This also results in Gaussian approximation.
- *Unscented Rauch-Tung-Striebel smoother (URTSS)* is the approximate non-linear smoothing algorithm corresponding to UKF.
- *Sequential Monte Carlo methods* or *particle filters and smoothers* represent the posterior distribution as a weighted set of Monte Carlo samples.

- *Unscented particle filter (UPF)* and *local linearization* based methods use UKFs and EKF, respectively, for approximating the importance distributions in sequential importance sampling.
- *Rao-Blackwellized particle filters and smoothers* use closed form integration (e.g., Kalman filters and RTS smoothers) for some of the state variables and Monte Carlo integration for others.
- *Interacting multiple models (IMM)*, and other *multiple model* methods approximate the posterior distributions with mixture Gaussian approximations.
- *Grid based methods* approximate the distribution as a discrete distribution defined in a finite grid.
- *Other methods* also exist, for example, based on series expansions, describing functions, basis function expansions, exponential family of distributions, variational Bayesian methods, batch Monte Carlo (e.g., MCMC), Galerkin approximations etc.

Chapter 2

From Bayesian Inference to Bayesian Optimal Filtering

2.1 Bayesian Inference

This section provides a brief presentation of the philosophical and mathematical foundations of Bayesian inference. The connections to the classical statistical inference are also briefly discussed.

2.1.1 Philosophy of Bayesian Inference

The purpose of Bayesian inference (Bernardo and Smith, 1994; Gelman et al., 1995) is to provide a mathematical machinery that can be used for modeling systems, where the uncertainties of the system are taken into account and the decisions are made according to rational principles. The tools of this machinery are the probability distributions and the rules of probability calculus.

If we compare the so called frequentist philosophy of statistical analysis to Bayesian inference the difference is that in Bayesian inference the probability of an event does not mean the proportion of the event in an infinite number of trials, but the uncertainty of the event in a single trial. Because models in Bayesian inference are formulated in terms of probability distributions, the probability axioms and computation rules of the probability theory (see, e.g., Shiryaev, 1996) also apply in the Bayesian inference.

2.1.2 Connection to Maximum Likelihood Estimation

Consider a situation, where we know the conditional distribution $p(\mathbf{y}_k | \boldsymbol{\theta})$ of conditionally independent random variables (measurements) $\mathbf{y}_1, \dots, \mathbf{y}_n$, but the parameter $\boldsymbol{\theta} \in \mathbb{R}^d$ is unknown. The classical statistical method for estimating the parameter is the *maximum likelihood method* (Milton and Arnold, 1995), where we maximize the joint probability of the measurements, also called the likelihood

function

$$L(\boldsymbol{\theta}) = \prod_k p(\mathbf{y}_k | \boldsymbol{\theta}). \quad (2.1)$$

The maximum of the likelihood function with respect to $\boldsymbol{\theta}$ gives the *maximum likelihood estimate* (ML-estimate)

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}). \quad (2.2)$$

The difference between the Bayesian inference and the maximum likelihood method is that the starting point of Bayesian inference is to formally consider the parameter $\boldsymbol{\theta}$ as a random variable. Then the posterior distribution of the parameter $\boldsymbol{\theta}$ can be computed by using the *Bayes' rule*

$$p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n) = \frac{p(\mathbf{y}_1, \dots, \mathbf{y}_n | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y}_1, \dots, \mathbf{y}_n)}, \quad (2.3)$$

where $p(\boldsymbol{\theta})$ is the prior distribution, which models the prior beliefs of the parameter before we have seen any data and $p(\mathbf{y}_1, \dots, \mathbf{y}_n)$ is a normalization term, which is independent of the parameter $\boldsymbol{\theta}$. Often this normalization constant is left out and if the measurements $\mathbf{y}_1, \dots, \mathbf{y}_n$ are conditionally independent given $\boldsymbol{\theta}$, the posterior distribution of the parameter can be written as

$$p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n) \propto p(\boldsymbol{\theta}) \prod_k p(\mathbf{y}_k | \boldsymbol{\theta}). \quad (2.4)$$

Because we are dealing with a distribution, we might now choose the most probable value of the random variable (MAP-estimate), which is given by the maximum of the posterior distribution. However, better estimate in mean squared sense is the posterior mean of the parameter (MMSE-estimate). There are an infinite number of other ways of choosing the point estimate from the distribution and the best way depends on the assumed loss function (or utility function). The ML-estimate can be considered as a MAP-estimate with uniform prior on the parameter $\boldsymbol{\theta}$.

One can also interpret Bayesian inference as a convenient method for including regularization terms into maximum likelihood estimation. The basic ML-framework does not have a self-consistent method for including regularization terms or prior information into statistical models. However, this regularization interpretation of Bayesian inference is not entirely right, because Bayesian inference is much more than this.

2.1.3 The Building Blocks of Bayesian Models

The basic blocks of a Bayesian model are the *prior model* containing the preliminary information on the parameter and the *likelihood model* determining the stochastic mapping from the parameter to the measurements. Using the combination rules, namely the Bayes' rule, it is possible to infer an estimate of the

parameters from the measurements. The distribution of the parameters, which is conditional to the observed measurements is called the *posterior distribution* and it is the distribution representing the state of knowledge about the parameters when all the information in the observed measurements and the model is used. *Predictive posterior distribution* is the distribution of the new (not yet observed) measurements when all the information in the observed measurements and the model is used.

- **Prior model**

The prior information consists of subjective experience based beliefs on the possible and impossible parameter values and their relative likelihoods before anything has been observed. The prior distribution is a mathematical representation of this information:

$$p(\boldsymbol{\theta}) = \text{Information on parameter } \boldsymbol{\theta} \text{ before seeing any observations.} \quad (2.5)$$

The lack of prior information can be expressed by using a non-informative prior. The non-informative prior distribution can be selected in various different ways (Gelman et al., 1995).

- **Likelihood model**

Between the true parameters and the measurements there often is a causal, but inaccurate or noisy relationship. This relationship is mathematically modeled using the likelihood distribution:

$$p(\mathbf{y} | \boldsymbol{\theta}) = \text{Distribution of observation } \mathbf{y} \text{ given the parameters } \boldsymbol{\theta}. \quad (2.6)$$

- **Posterior distribution**

Posterior distribution is the conditional distribution of the parameters, and it represents the information we have after the measurement \mathbf{y} has been obtained. It can be computed by using the Bayes' rule:

$$p(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{y})} \propto p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta}), \quad (2.7)$$

where the normalization constant is given as

$$p(\mathbf{y}) = \int_{\mathbb{R}^d} p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (2.8)$$

In the case of multiple measurements $\mathbf{y}_1, \dots, \mathbf{y}_n$, if the measurements are conditionally independent the joint likelihood of all measurements is the product of individual measurements and the posterior distribution is

$$p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n) \propto p(\boldsymbol{\theta}) \prod_k p(\mathbf{y}_k | \boldsymbol{\theta}), \quad (2.9)$$

where the normalization term can be computed by integrating the right hand side over $\boldsymbol{\theta}$. If the random variable is discrete the integration reduces to summation.

- **Predictive posterior distribution**

The predictive posterior distribution is the distribution of new measurements \mathbf{y}_{n+1} :

$$p(\mathbf{y}_{n+1} | \mathbf{y}_1, \dots, \mathbf{y}_n) = \int_{\mathbb{R}^d} p(\mathbf{y}_{n+1} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n) d\boldsymbol{\theta}. \quad (2.10)$$

After obtaining the measurements $\mathbf{y}_1, \dots, \mathbf{y}_n$ the predictive posterior distribution can be used for computing the probability distribution for $n + 1$:th measurement, which has not been observed yet.

In the case of tracking, we could imagine that the parameter is the sequence of dynamic states of a target, where the state contains the position and velocity. Or in the continuous-discrete setting the parameter would be an infinite-dimensional random function describing the trajectory of the target at a given time interval. In both cases the measurements could be, for example, noisy distance and direction measurements produced by a radar.

2.1.4 Bayesian Point Estimates

The distributions as such have no use in applications, but also in Bayesian computations finite dimensional summaries (point estimates) are needed. This selection of a point from space based on observed values of random variables is a statistical decision, and therefore this selection procedure is most naturally formulated in terms of *statistical decision theory* (Berger, 1985; Bernardo and Smith, 1994; Raiffa and Schlaifer, 2000).

Definition 2.1 (Loss Function). *A loss function $L(\boldsymbol{\theta}, \mathbf{a})$ is a scalar valued function, which determines the loss of taking the action \mathbf{a} , when the true parameter value is $\boldsymbol{\theta}$. The action (or control) is the statistical decision to be made based on the currently available information.*

Instead of loss functions it is also possible to work with utility functions $U(\boldsymbol{\theta}, \mathbf{a})$, which determine the reward from taking the action \mathbf{a} with parameter values $\boldsymbol{\theta}$. Loss functions can be converted to utility functions and vice versa by defining $U(\boldsymbol{\theta}, \mathbf{a}) = -L(\boldsymbol{\theta}, \mathbf{a})$.

If the value of parameter $\boldsymbol{\theta}$ is not known, but the knowledge on the parameter can be expressed in terms of the posterior distribution $p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n)$, then the natural choice is the action, which gives the *minimum (maximum) of the expected loss (utility)* (Berger, 1985):

$$\mathbb{E}[L(\boldsymbol{\theta}, \mathbf{a}) | \mathbf{y}_1, \dots, \mathbf{y}_n] = \int_{\mathbb{R}^d} L(\boldsymbol{\theta}, \mathbf{a}) p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n) d\boldsymbol{\theta}. \quad (2.11)$$

Commonly used loss functions are the following:

- *Quadratic error loss*: If the loss function is quadratic

$$L(\boldsymbol{\theta}, \mathbf{a}) = (\boldsymbol{\theta} - \mathbf{a})^T (\boldsymbol{\theta} - \mathbf{a}), \quad (2.12)$$

then the optimal choice \mathbf{a}_o is the *posterior mean* of the distribution of $\boldsymbol{\theta}$:

$$\mathbf{a}_o = \int_{\mathbb{R}^d} \boldsymbol{\theta} p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n) d\boldsymbol{\theta}. \quad (2.13)$$

This posterior mean based estimate is often called the *minimum mean squared error (MMSE)* estimate of the parameter $\boldsymbol{\theta}$. The quadratic loss is the most commonly used loss function, because it is easy to handle mathematically and because in the case of Gaussian posterior distribution the MAP estimate and the median coincide with the posterior mean.

- *Absolute error loss*: The loss function of the form

$$L(\boldsymbol{\theta}, \mathbf{a}) = \sum_i |\theta_i - a_i|, \quad (2.14)$$

is called an absolute error loss and in this case the optimal choice is the *median* of the distribution (i.e., medians of the marginal distributions in multidimensional case).

- *0-1 loss*: If the loss function is of the form

$$L(\boldsymbol{\theta}, \mathbf{a}) = \begin{cases} 1 & , \text{ if } \boldsymbol{\theta} = \mathbf{a} \\ 0 & , \text{ if } \boldsymbol{\theta} \neq \mathbf{a} \end{cases} \quad (2.15)$$

then the optimal choice is the maximum of the posterior distribution, that is, the *maximum a posteriori (MAP)* estimate of the parameter.

2.1.5 Numerical Methods

In principle, Bayesian inference provides the equations for computing the posterior distributions and point estimates for any model once the model specification has been set up. However, the practical problem is that computation of the integrals involved in the equations can rarely be performed analytically and numerical methods are needed. Here we shall briefly describe numerical methods, which are also applicable in higher dimensional problems: Gaussian approximations, multi-dimensional quadratures, Monte Carlo methods, and importance sampling.

- Very common types of approximations are *Gaussian approximations* (Gelman et al., 1995), where the posterior distribution is approximated with a Gaussian distribution

$$p(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n) \approx \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}, \mathbf{P}). \quad (2.16)$$

The mean \mathbf{m} and covariance \mathbf{P} of the Gaussian approximation can be either computed by matching the first two moments of the posterior distribution, or by using the maximum of the distribution as the mean estimate and approximating the covariance with the curvature of the posterior on the mode.

- *Multi-dimensional quadrature or cubature methods* such as Gauss-Hermite quadrature can also be often used if the dimensionality of the integral is moderate. In those methods the idea is to deterministically form a representative set of sample points $\Theta = \{\boldsymbol{\theta}^{(i)} \mid i = 1, \dots, N\}$ (sometimes called *sigma points*) and form the approximation of the integral as weighted average:

$$E[\mathbf{g}(\boldsymbol{\theta}) \mid \mathbf{y}_1, \dots, \mathbf{y}_n] \approx \sum_{i=1}^N W^{(i)} \mathbf{g}(\boldsymbol{\theta}^{(i)}), \quad (2.17)$$

where the numerical values of the weights $W^{(i)}$ are determined by the algorithm. The sample points and weights can be selected, for example, to give exact answers for polynomials up to certain degree or to account for the moments up to certain degree.

- In direct *Monte Carlo methods* a set of N samples from the posterior distribution is randomly drawn

$$\boldsymbol{\theta}^{(i)} \sim p(\boldsymbol{\theta} \mid \mathbf{y}_1, \dots, \mathbf{y}_n), \quad i = 1, \dots, N, \quad (2.18)$$

and expectation of any function $\mathbf{g}(\cdot)$ can be then approximated as the sample average

$$E[\mathbf{g}(\boldsymbol{\theta}) \mid \mathbf{y}_1, \dots, \mathbf{y}_n] \approx \frac{1}{N} \sum_i \mathbf{g}(\boldsymbol{\theta}^{(i)}). \quad (2.19)$$

Another interpretation of this is that Monte Carlo methods form an approximation of the posterior density of the form

$$p(\boldsymbol{\theta} \mid \mathbf{y}_1, \dots, \mathbf{y}_n) \approx \frac{1}{N} \sum_{i=1}^N \delta(x - x^{(i)}), \quad (2.20)$$

where $\delta(\cdot)$ is the Dirac delta function. The convergence of Monte Carlo approximation is guaranteed by the *central limit theorem (CLT)* (see, e.g., Liu, 2001) and the error term is, at least in theory, independent of the dimensionality of $\boldsymbol{\theta}$.

- Efficient methods for generating non-independent Monte Carlo samples are the *Markov chain Monte Carlo (MCMC)* methods (see, e.g., Gilks et al., 1996). In MCMC methods, a Markov chain is constructed such that it has the target distribution as its stationary distribution. By simulating the Markov chain, samples from the target distribution can be generated.
- *Importance sampling* (see, e.g., Liu, 2001) is a simple algorithm for generating *weighted* samples from the target distribution. The difference to the direct Monte Carlo sampling and to MCMC is that each of the particles contains a weight, which corrects the difference between the actual target

distribution and the approximation obtained from an importance distribution $\pi(\cdot)$.

Importance sampling estimate can be formed by drawing N samples from the *importance distribution*

$$\boldsymbol{\theta}^{(i)} \sim \pi(\boldsymbol{\theta} | \mathbf{y}_1, \dots, \mathbf{y}_n), \quad i = 1, \dots, N. \quad (2.21)$$

The *importance weights* are then computed as

$$w^{(i)} = \frac{p(\boldsymbol{\theta}^{(i)} | \mathbf{y}_1, \dots, \mathbf{y}_n)}{\pi(\boldsymbol{\theta}^{(i)} | \mathbf{y}_1, \dots, \mathbf{y}_n)}, \quad (2.22)$$

and the expectation of any function $\mathbf{g}(\cdot)$ can be then approximated as

$$\mathbb{E}[\mathbf{g}(\boldsymbol{\theta}) | \mathbf{y}_1, \dots, \mathbf{y}_n] \approx \frac{\sum_{i=1}^N w^{(i)} \mathbf{g}(\boldsymbol{\theta}^{(i)})}{\sum_{i=1}^N w^{(i)}}. \quad (2.23)$$

2.2 Batch and Recursive Estimation

In order to understand the meaning and applicability of optimal filtering and its relationship with recursive estimation, it is useful to go through an example, where we solve a simple and familiar linear regression problem in a recursive manner. After that we shall generalize this concept to include a dynamic model in order to illustrate the differences in dynamic and batch estimation.

2.2.1 Batch Linear Regression

Consider the linear regression model

$$y_k = \theta_1 + \theta_2 t_k + \epsilon_k, \quad (2.24)$$

where we assume that the measurement noise is zero mean Gaussian with a given variance $\epsilon_k \sim \mathcal{N}(0, \sigma^2)$ and the prior distribution for parameters is Gaussian with known mean and covariance $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{m}_0, \mathbf{P}_0)$. In the classical linear regression problem we want to estimate the parameters $\boldsymbol{\theta} = (\theta_1 \ \theta_2)^T$ from a set of measurement data $\mathcal{D} = \{(y_1, t_1), \dots, (y_K, t_K)\}$. The measurement data and the true linear function used in simulation are illustrated in Figure 2.1.

In compact probabilistic notation the linear regression model can be written as

$$\begin{aligned} p(y_k | \boldsymbol{\theta}) &= \mathcal{N}(y_k | \mathbf{H}_k \boldsymbol{\theta}, \sigma^2) \\ p(\boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{P}_0). \end{aligned} \quad (2.25)$$

where we have introduced the matrix $\mathbf{H}_k = (1 \ t_k)$ and $\mathcal{N}(\cdot)$ denotes the Gaussian probability density function (see, Appendix A.1). The likelihood of y_k is, of course, conditional on the regressors t_k also (or equivalently \mathbf{H}_k), but we will not

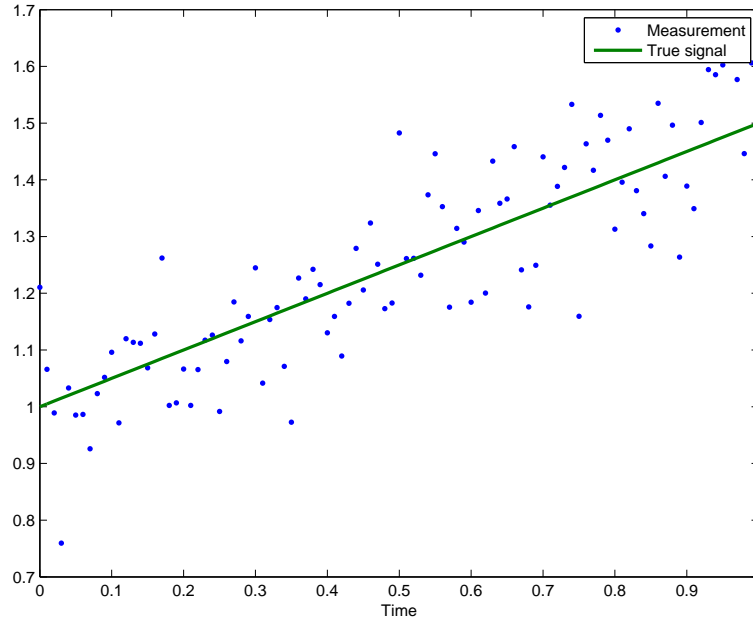


Figure 2.1: The underlying truth and the measurement data in the simple linear regression problem.

denote this dependence explicitly to simplify the notation and from now on this dependence is assumed to be understood from the context.

The *batch solution* to this linear regression problem can be obtained by straightforward application of the Bayes' rule:

$$\begin{aligned} p(\boldsymbol{\theta} | y_{1:k}) &\propto p(\boldsymbol{\theta}) \prod_k p(y_k | \boldsymbol{\theta}) \\ &= \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_0, \mathbf{P}_0) \prod_k \mathcal{N}(y_k | \mathbf{H}_k \boldsymbol{\theta}, \sigma^2). \end{aligned}$$

Also in the posterior distribution above, we assume the conditioning on t_k and \mathbf{H}_k , but will not denote it explicitly. Thus the posterior distribution is denoted to be conditional on $y_{1:k} = \{y_1, \dots, y_k\}$, and not on the data set \mathcal{D} containing the regressor values t_k also. The reason for this simplification is that the simplified notation will also work in more general filtering problems, where there is no natural way of defining the associated regressor variables.

Because the prior and likelihood are Gaussian, the posterior distribution will also be Gaussian:

$$p(\boldsymbol{\theta} | y_{1:k}) = \mathcal{N}(\boldsymbol{\theta} | \mathbf{m}_K, \mathbf{P}_K). \quad (2.26)$$

The mean and covariance can be obtained by completing the quadratic form in the

exponent, which gives:

$$\begin{aligned} \mathbf{m}_K &= \left[\mathbf{P}_0^{-1} + \frac{1}{\sigma^2} \mathbf{H}^T \mathbf{H} \right]^{-1} \left[\frac{1}{\sigma^2} \mathbf{H}^T \mathbf{y} + \mathbf{P}_0^{-1} \mathbf{m}_0 \right] \\ \mathbf{P}_K &= \left[\mathbf{P}_0^{-1} + \frac{1}{\sigma^2} \mathbf{H}^T \mathbf{H} \right]^{-1}, \end{aligned} \quad (2.27)$$

where $\mathbf{H}_k = (1 \ t_k)$ and

$$\mathbf{H} = \begin{pmatrix} \mathbf{H}_1 \\ \vdots \\ \mathbf{H}_K \end{pmatrix} = \begin{pmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_K \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_K \end{pmatrix}. \quad (2.28)$$

Figure 2.2 shows the result of batch linear regression, where the posterior mean parameter values are used as the linear regression parameters.

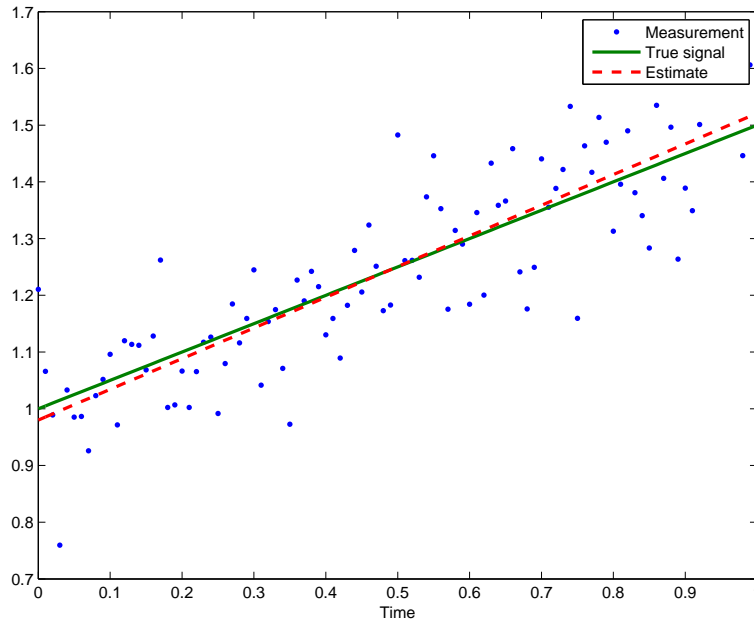


Figure 2.2: The result of simple linear regression with a slight regularization prior used for the regression parameters. For simplicity, the variance was assumed to be known.

2.2.2 Recursive Linear Regression

Recursive solution to the regression problem (2.25) can be obtained by assuming that we already have obtained posterior distribution conditioned on the previous measurements $1, \dots, k-1$:

$$p(\boldsymbol{\theta} | y_{1:k-1}) = N(\boldsymbol{\theta} | \mathbf{m}_{k-1}, \mathbf{P}_{k-1}).$$

Now assume that we have obtained a new measurement y_k and we want to compute the posterior distribution of $\boldsymbol{\theta}$ given the old measurements $y_{1:k-1}$ and the new measurement y_k . According to the model specification the new measurement has the likelihood

$$p(y_k | \boldsymbol{\theta}) = \text{N}(y_k | \mathbf{H}_k \boldsymbol{\theta}, \sigma^2).$$

Using the batch version equations such that we interpret the previous posterior as the prior, we can calculate the distribution

$$\begin{aligned} p(\boldsymbol{\theta} | y_{1:k}) &\propto p(y_k | \boldsymbol{\theta}) p(\boldsymbol{\theta} | y_{1:k-1}) \\ &\propto \text{N}(\boldsymbol{\theta} | \mathbf{m}_k, \mathbf{P}_k), \end{aligned} \quad (2.29)$$

where the Gaussian distribution parameters are

$$\begin{aligned} \mathbf{m}_k &= \left[\mathbf{P}_{k-1}^{-1} + \frac{1}{\sigma^2} \mathbf{H}_k^T \mathbf{H}_k \right]^{-1} \left[\frac{1}{\sigma^2} \mathbf{H}_k^T y_k + \mathbf{P}_{k-1}^{-1} \mathbf{m}_{k-1} \right] \\ \mathbf{P}_k &= \left[\mathbf{P}_{k-1}^{-1} + \frac{1}{\sigma^2} \mathbf{H}_k^T \mathbf{H}_k \right]^{-1}. \end{aligned} \quad (2.30)$$

By using the matrix inversion lemma, the covariance calculation can be written as

$$\mathbf{P}_k = \mathbf{P}_{k-1} - \mathbf{P}_{k-1} \mathbf{H}_k^T [\mathbf{H}_k \mathbf{P}_{k-1} \mathbf{H}_k^T + \sigma^2]^{-1} \mathbf{H}_k \mathbf{P}_{k-1}.$$

By introducing temporary variables S_k and \mathbf{K}_k the calculation of mean and covariance can be written in form

$$\begin{aligned} S_k &= \mathbf{H}_k \mathbf{P}_{k-1} \mathbf{H}_k^T + \sigma^2 \\ \mathbf{K}_k &= \mathbf{P}_{k-1} \mathbf{H}_k^T S_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_{k-1} + \mathbf{K}_k [y_k - \mathbf{H}_k \mathbf{m}_{k-1}] \\ \mathbf{P}_k &= \mathbf{P}_{k-1} - \mathbf{K}_k S_k \mathbf{K}_k^T. \end{aligned} \quad (2.31)$$

Note that $S_k = \mathbf{H}_k \mathbf{P}_{k-1} \mathbf{H}_k^T + \sigma^2$ is a scalar, because measurements are scalar and thus no matrix inversion is required.

The equations above actually are special cases of the Kalman filter update equations. Only the update part of the equations is required, because the estimated parameters are assumed to be constant, that is, there is no a priori stochastic dynamics model for the parameters $\boldsymbol{\theta}$. Figure 2.3 illustrates the convergence of the means and variances of parameters during the recursive estimation.

2.2.3 Batch vs. Recursive Estimation

In this section we shall generalize the recursion idea used in the previous section to general probabilistic models. The underlying idea is simply that at each measurement we treat the posterior distribution of previous time step as the prior for the current time step. This way we can compute the same solution in recursive manner

that we would obtain by direct application of Bayesian rule to the whole (batch) data set.

The *batch Bayesian solution* to a statistical estimation problem can be formulated as follows:

1. Specify the likelihood model of measurements $p(\mathbf{y}_k | \boldsymbol{\theta})$ given the parameter $\boldsymbol{\theta}$. Typically the measurements \mathbf{y}_k are assumed to be conditionally independent such that

$$p(\mathbf{y}_{1:K} | \boldsymbol{\theta}) = \prod_k p(\mathbf{y}_k | \boldsymbol{\theta}).$$

2. The prior information about the parameter $\boldsymbol{\theta}$ is encoded into the prior distribution $p(\boldsymbol{\theta})$.
3. The observed data set is $\mathcal{D} = \{(t_1, \mathbf{y}_1), \dots, (t_K, \mathbf{y}_K)\}$, or if we drop the explicit conditioning to t_k , the data is $\mathcal{D} = \mathbf{y}_{1:K}$.
4. The batch Bayesian solution to the statistical estimation problem can be computed by applying the Bayes' rule

$$p(\boldsymbol{\theta} | \mathbf{y}_{1:K}) = \frac{1}{Z} p(\boldsymbol{\theta}) \prod_k p(\mathbf{y}_k | \boldsymbol{\theta}).$$

For example, the batch solution of the above kind to the linear regression problem (2.25) was given by Equations (2.26) and (2.27).

The *recursive Bayesian solution* to the above statistical estimation problem can be formulated as follows:

1. The distribution of measurements is again modeled by the likelihood function $p(\mathbf{y}_k | \boldsymbol{\theta})$ and the measurements are assumed to be conditionally independent.
2. In the beginning of estimation (i.e, at step 0), all the information about the parameter $\boldsymbol{\theta}$ we have, is the prior distribution $p(\boldsymbol{\theta})$.
3. The measurements are assumed to be obtained one at a time, first \mathbf{y}_1 , then \mathbf{y}_2 and so on. At each step we use the posterior distribution from the previous time step as the current prior distribution:

$$\begin{aligned} p(\boldsymbol{\theta} | \mathbf{y}_1) &= \frac{1}{Z_1} p(\mathbf{y}_1 | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \\ p(\boldsymbol{\theta} | \mathbf{y}_{1:2}) &= \frac{1}{Z_2} p(\mathbf{y}_2 | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}_1) \\ p(\boldsymbol{\theta} | \mathbf{y}_{1:3}) &= \frac{1}{Z_3} p(\mathbf{y}_3 | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}_{1:2}) \\ &\vdots \\ p(\boldsymbol{\theta} | \mathbf{y}_{1:K}) &= \frac{1}{Z_K} p(\mathbf{y}_K | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}_{1:K-1}). \end{aligned}$$

It is easy to show that the posterior distribution at the final step above is exactly the posterior distribution obtained by the batch solution. Also, re-ordering of measurements does not change the final solution.

For example, the Equations (2.29) and (2.30) give the one step update rule for the linear regression problem in Equation (2.25).

The recursive formulation of Bayesian estimation has many useful properties:

- The recursive solution can be considered as the *online learning* solution to the Bayesian learning problem. That is, the information on the parameters is updated in online manner using new pieces of information as they arrive.
- Because each step in the recursive estimation is a full Bayesian update step, *batch* Bayesian inference is a *special case of recursive* Bayesian inference.
- Due to the sequential nature of estimation we can also model the effect of time to parameters. That is, we can build model to what happens to the parameter θ between the measurements – this is actually the *basis of filtering theory*, where time behavior is modeled by assuming the parameter to be a time-dependent stochastic process $\theta(t)$.

2.3 Towards Bayesian Filtering

Now that we are able to solve the static linear regression problem in recursive manner, we can proceed towards Bayesian filtering by allowing the parameters change between the measurements. By generalizing this idea, we encounter the Kalman filter, which is the workhorse of dynamic estimation.

2.3.1 Drift Model for Linear Regression

Assume that we have similar linear regression model as in Equation (2.25), but the parameter θ is allowed to perform *Gaussian random walk* between the measurements:

$$\begin{aligned} p(y_k | \theta_k) &= N(y_k | \mathbf{H}_k \theta_k, \sigma^2) \\ p(\theta_k | \theta_{k-1}) &= N(\theta_k | \theta_{k-1}, \mathbf{Q}) \\ p(\theta_0) &= N(\theta_0 | \mathbf{m}_0, \mathbf{P}_0), \end{aligned} \tag{2.32}$$

where \mathbf{Q} is the covariance of the random walk. Now, given the distribution

$$p(\theta_{k-1} | y_{1:k-1}) = N(\theta_{k-1} | \mathbf{m}_{k-1}, \mathbf{P}_{k-1}),$$

the joint distribution of θ_k and θ_{k-1} is¹

$$p(\theta_k, \theta_{k-1} | y_{1:k-1}) = p(\theta_k | \theta_{k-1}) p(\theta_{k-1} | y_{1:k-1}).$$

¹Note that this formula is correct only for Markovian dynamic models, where $p(\theta_k | \theta_{k-1}, y_{1:k-1}) = p(\theta_k | \theta_{k-1})$.

The distribution of $\boldsymbol{\theta}_k$ given the measurement history up to time step $k - 1$ can be calculated by integrating over $\boldsymbol{\theta}_{k-1}$

$$p(\boldsymbol{\theta}_k | y_{1:k-1}) = \int p(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1}) p(\boldsymbol{\theta}_{k-1} | y_{1:k-1}) d\boldsymbol{\theta}_{k-1}.$$

This relationship is sometimes called the *Chapman-Kolmogorov equation*. Because $p(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1})$ and $p(\boldsymbol{\theta}_{k-1} | y_{1:k-1})$ are Gaussian, the result of the marginalization is Gaussian:

$$p(\boldsymbol{\theta}_k | y_{1:k-1}) = \mathcal{N}(\boldsymbol{\theta}_k | \mathbf{m}_k^-, \mathbf{P}_k^-),$$

where

$$\begin{aligned} \mathbf{m}_k^- &= \mathbf{m}_{k-1} \\ \mathbf{P}_k^- &= \mathbf{P}_{k-1} + \mathbf{Q}. \end{aligned}$$

By using this as the prior distribution for the measurement likelihood $p(y_k | \boldsymbol{\theta}_k)$ we get the parameters of the posterior distribution

$$p(\boldsymbol{\theta}_k | y_{1:k}) = \mathcal{N}(\boldsymbol{\theta}_k | \mathbf{m}_k, \mathbf{P}_k),$$

which are given by equations (2.31), when \mathbf{m}_{k-1} and \mathbf{P}_{k-1} are replaced by \mathbf{m}_k^- and \mathbf{P}_k^- :

$$\begin{aligned} S_k &= \mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \sigma^2 \\ \mathbf{K}_k &= \mathbf{P}_k^- \mathbf{H}_k^T S_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k [y_k - \mathbf{H}_k \mathbf{m}_k^-] \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k S_k \mathbf{K}_k^T. \end{aligned} \tag{2.33}$$

This recursive computational algorithm for the time-varying linear regression weights is again a special case of the Kalman filter algorithm. Figure 2.4 shows the result of recursive estimation of sine signal assuming a small diagonal Gaussian drift model for the parameters.

At this point we shall change from the *regression notation* used so far into *state space model notation*, which is commonly used in Kalman filtering and related dynamic estimation literature. Because this notation easily causes confusion to people who have got used to regression notation, this point is emphasized:

- In *state space notation* \mathbf{x} means the unknown state of the system, that is, the vector of *unknown parameters in the system*. It is *not* the regressor, covariate or input variable of the system.
- For example, the time-varying linear regression model with drift presented in this section can be transformed into more standard *state space model notation* by replacing the variable $\boldsymbol{\theta}_k = (\theta_{1,k} \ \theta_{2,k})^T$ with the variable $\mathbf{x}_k = (x_{1,k} \ x_{2,k})^T$:

$$\begin{aligned} p(y_k | \mathbf{x}_k) &= \mathcal{N}(y_k | \mathbf{H}_k \mathbf{x}_k, \sigma^2) \\ p(\mathbf{x}_k | \mathbf{x}_{k-1}) &= \mathcal{N}(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{Q}) \\ p(\mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_0 | \mathbf{m}_0, \mathbf{P}_0). \end{aligned} \tag{2.34}$$

2.3.2 Kalman Filter

The linear model with drift in the previous section had the disadvantage that the covariates t_k occurred explicitly in the model specification. The problem with this is that when we get more and more measurements, the parameter t_k grows without a bound. Thus the conditioning of the problem also gets worse in time. For practical reasons it also would be desirable to have time-invariant model, that is, a model which is not dependent on the absolute time, but only on the relative positions of states and measurements in time.

The alternative state space formulation of the linear model with drift, without using explicit covariates can be done as follows. Let's denote time difference between consecutive times as $\Delta t_{k-1} = t_k - t_{k-1}$. The idea is that if the underlying phenomenon (signal, state, parameter) x_k was exactly linear, the difference between adjacent time points could be written exactly as

$$x_k - x_{k-1} = \dot{x} \Delta t_{k-1} \quad (2.35)$$

where \dot{x} is the derivative, which is constant in the exactly linear case. The divergence from the exactly linear function can be modeled by assuming that the above equation does not hold exactly, but there is a small noise term on the right hand side. The derivative can also be assumed to perform small random walk and thus not be exactly constant. This model can be written as follows:

$$\begin{aligned} x_{1,k} &= x_{1,k-1} + \Delta t_{k-1} \dot{x}_{2,k-1} + w_1 \\ x_{2,k} &= x_{2,k-1} + w_2 \\ y_k &= x_{1,k} + e, \end{aligned} \quad (2.36)$$

where the signal is the first components of the state $x_{1,k}$ and the derivative is the second $x_{2,k}$. The noises are $e \sim N(0, \sigma^2)$, $(w_1; w_2) \sim N(0, \mathbf{Q})$. The model can also be written in form

$$\begin{aligned} p(y_k | \mathbf{x}_k) &= N(y_k | \mathbf{H} \mathbf{x}_k, \sigma^2) \\ p(\mathbf{x}_k | \mathbf{x}_{k-1}) &= N(\mathbf{x}_k | \mathbf{A}_{k-1} \mathbf{x}_{k-1}, \mathbf{Q}), \end{aligned} \quad (2.37)$$

where

$$\mathbf{A}_{k-1} = \begin{pmatrix} 1 & \Delta t_{k-1} \\ 0 & 1 \end{pmatrix}, \quad \mathbf{H} = (1 \ 0).$$

With suitable \mathbf{Q} this model is actually equivalent to model (2.32), but in this formulation we explicitly estimate the state of the signal (point on the regression line) instead of the linear regression parameters.

We could now explicitly derive the recursion equations in the same manner as we did in the previous sections. However, we can also use the *Kalman filter*, which is a readily derived recursive solution to generic linear Gaussian models of the form

$$\begin{aligned} p(\mathbf{y}_k | \mathbf{x}_k) &= N(\mathbf{y}_k | \mathbf{H}_k \mathbf{x}_k, \mathbf{R}_k) \\ p(\mathbf{x}_k | \mathbf{x}_{k-1}) &= N(\mathbf{x}_k | \mathbf{A}_{k-1} \mathbf{x}_{k-1}, \mathbf{Q}_{k-1}). \end{aligned}$$

Our alternative linear regression model in Equation (2.36) can be seen to be a special case of these models. The Kalman filter equations are often expressed as prediction and update steps as follows:

1. *Prediction step:*

$$\begin{aligned}\mathbf{m}_k^- &= \mathbf{A}_{k-1} \mathbf{m}_{k-1} \\ \mathbf{P}_k^- &= \mathbf{A}_{k-1} \mathbf{P}_{k-1} \mathbf{A}_{k-1}^T + \mathbf{Q}_{k-1}.\end{aligned}$$

2. *Update step:*

$$\begin{aligned}\mathbf{S}_k &= \mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k \\ \mathbf{K}_k &= \mathbf{P}_k^- \mathbf{H}_k^T \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k [\mathbf{y}_k - \mathbf{H}_k \mathbf{m}_k^-] \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.\end{aligned}$$

The result of tracking the sine signal with Kalman filter is shown in Figure 2.5. All the mean and covariance calculation equations given in this document so far have been special cases of the above equations, including the batch solution to the scalar measurement case (which is a one-step solution). The Kalman filter recursively computes the mean and covariance of the posterior distributions of the form

$$p(\mathbf{x}_k | \mathbf{y}_1, \dots, \mathbf{y}_k) = \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k).$$

Note that the estimates of \mathbf{x}_k derived from this distribution are non-anticipative in the sense that they are only conditional to measurements obtained before and at the time step k . However, after we have obtained measurements $\mathbf{y}_1, \dots, \mathbf{y}_k$, we could compute estimates of $\mathbf{x}_{k-1}, \mathbf{x}_{k-2}, \dots$, which are also conditional to the measurements after the corresponding state time steps. Because more measurements and more information is available for the estimator, these estimates can be expected to be more accurate than the non-anticipative measurements computed by the filter.

The above mentioned problem of computing estimates of state by conditioning not only to previous measurements, but also to future measurements is called *optimal smoothing* as already mentioned in Section 1.2.3. The optimal smoothing solution to the linear Gaussian state space models is given by the *Rauch-Tung-Striebel smoother*. The full Bayesian theory of optimal smoothing as well as the related algorithms will be presented in Chapter 4.

It is also possible to predict the time behavior of the state in the future that we have not yet measured. This procedure is called *optimal prediction*. Because optimal prediction can always be done by iterating the prediction step of the optimal filter, no specialized algorithms are needed for this.

The non-linear generalizations of optimal prediction, filtering and smoothing can be obtained by replacing the Gaussian distributions and linear functions in

model (2.37) with non-Gaussian and non-linear ones. The Bayesian dynamic estimation theory described in this document can be applied to generic non-linear filtering models of the following form:

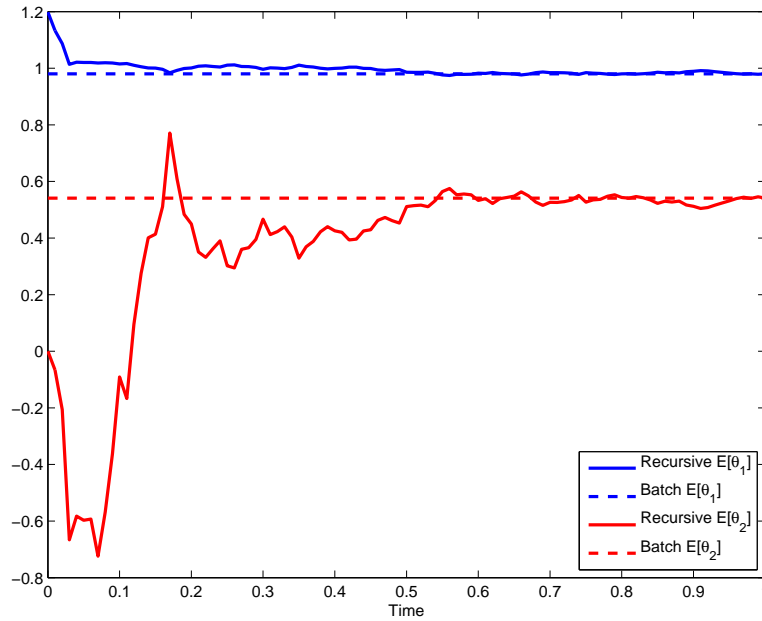
$$\begin{aligned}\mathbf{y}_k &\sim p(\mathbf{y}_k | \mathbf{x}_k) \\ \mathbf{x}_k &\sim p(\mathbf{x}_k | \mathbf{x}_{k-1}).\end{aligned}$$

To understand the generality of this model is it useful to note that if we dropped the time-dependence from the state we would get the model

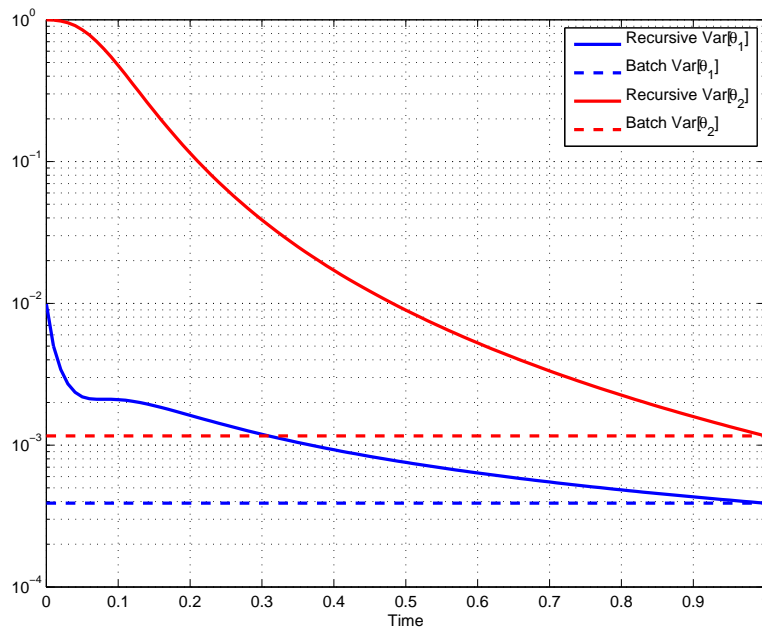
$$\begin{aligned}\mathbf{y}_k &\sim p(\mathbf{y}_k | \mathbf{x}) \\ \mathbf{x} &\sim p(\mathbf{x}).\end{aligned}$$

Because \mathbf{x} denotes an arbitrary set of parameters or hyper-parameters of the system, all static Bayesian models are special cases of this model. Thus in dynamic estimation context we extend the static models by allowing a Markov model for the time-behavior of the (hyper)parameters.

The Markovianity also is less of a restriction than it sounds, because what we have is a vector valued Markov process, not a scalar one. The reader may recall from the elementary calculus that differential equations of an arbitrary order can be always transformed into vector valued differential equations of the first order. In analogous manner, Markov processes of an arbitrary order can be transformed into vector valued first order Markov processes.



(a)



(b)

Figure 2.3: (a) Convergence of the recursive linear regression means. The final value is exactly the same as that was obtained with batch linear regression. Note that time has been scaled to 1 at $k = K$. (b) Convergence of the variances plotted on logarithmic scale. As can be seen, every measurement brings more information and the uncertainty decreases monotonically. The final values are the same as the variances obtained from the batch solution.

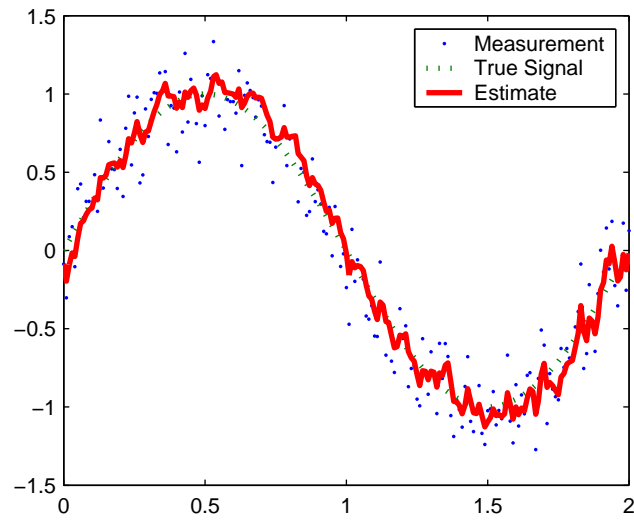


Figure 2.4: Example of tracking sine signal with linear model with drift, where the parameters are allowed to vary according to Gaussian random walk model.

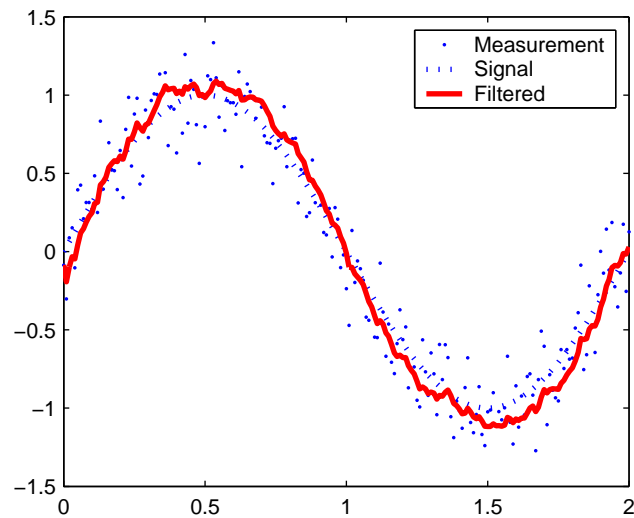


Figure 2.5: Example of tracking sine signal with locally linear state space model. The result differs a bit from the random walk parameter model, because of slightly different choice of process noise. It could be made equivalent if desired.

Chapter 3

Optimal Filtering

In this chapter we first present the classical formulation of the discrete-time optimal filtering as recursive Bayesian inference. Then the classical Kalman filters, extended Kalman filters and statistical linearization based filters are presented in terms of the general theory. In addition to the classical algorithms the unscented Kalman filter and general assumed density filters are also presented. Sequential importance resampling based particle filtering, as well as Rao-Blackwellized particle filtering are also covered.

3.1 Formal Filtering Equations and Exact Solutions

3.1.1 Discrete-Time Probabilistic State Space Models

Before going into the practical non-linear filtering algorithms, in the next sections the theory of probabilistic (Bayesian) filtering is presented. The Kalman filtering equations, which are the closed form solutions to the linear Gaussian discrete-time optimal filtering problem, are also derived.

Definition 3.1 (Discrete-time state space model). Discrete-time state space model is a recursively defined probabilistic model of the form

$$\begin{aligned}\mathbf{x}_k &\sim p(\mathbf{x}_k | \mathbf{x}_{k-1}) \\ \mathbf{y}_k &\sim p(\mathbf{y}_k | \mathbf{x}_k),\end{aligned}\tag{3.1}$$

where

- $\mathbf{x}_k \in \mathbb{R}^n$ is the state of the system on the time step k .
- $\mathbf{y}_k \in \mathbb{R}^m$ is the measurement on the time step k .
- $p(\mathbf{x}_k | \mathbf{x}_{k-1})$ is the dynamic model, which models the stochastic dynamics of the system. The dynamic model can be a probability density, a counting measure or combination of them depending on if the state \mathbf{x}_k is continuous, discrete or hybrid.

- $p(\mathbf{y}_k | \mathbf{x}_k)$ is the measurement model, which models the distribution of the measurements given the state.

The model is assumed to be Markovian, which means that it has the following two properties:

Property 3.1 (Markov property of states).

States $\{\mathbf{x}_k : k = 1, 2, \dots\}$ form a Markov sequence (or Markov chain if the state is discrete). This Markov property means that \mathbf{x}_k (and actually the whole future $\mathbf{x}_{k+1}, \mathbf{x}_{k+2}, \dots$) given \mathbf{x}_{k-1} is independent from anything that has happened in the past:

$$p(\mathbf{x}_k | \mathbf{x}_{1:k-1}, \mathbf{y}_{1:k-1}) = p(\mathbf{x}_k | \mathbf{x}_{k-1}). \quad (3.2)$$

Also the past is independent of the future given the present:

$$p(\mathbf{x}_{k-1} | \mathbf{x}_{k:T}, \mathbf{y}_{k:T}) = p(\mathbf{x}_{k-1} | \mathbf{x}_k). \quad (3.3)$$

Property 3.2 (Conditional independence of measurements).

The measurement \mathbf{y}_k given the state \mathbf{x}_k is conditionally independent from the measurement and state histories:

$$p(\mathbf{y}_k | \mathbf{x}_{1:k}, \mathbf{y}_{1:k-1}) = p(\mathbf{y}_k | \mathbf{x}_k). \quad (3.4)$$

As simple example of a Markovian sequence is the Gaussian random walk. When this is combined with noisy measurements, we obtain an example of a probabilistic state space model:

Example 3.1 (Gaussian random walk). *Gaussian random walk model can be written as*

$$\begin{aligned} x_k &= x_{k-1} + w_{k-1}, & w_{k-1} &\sim \mathcal{N}(0, q) \\ y_k &= x_k + e_k, & e_k &\sim \mathcal{N}(0, r), \end{aligned} \quad (3.5)$$

where x_k is the hidden state and y_k is the measurement. In terms of probability densities the model can be written as

$$\begin{aligned} p(x_k | x_{k-1}) &= \mathcal{N}(x_k | x_{k-1}, q) \\ &= \frac{1}{\sqrt{2\pi q}} \exp\left(-\frac{1}{2q}(x_k - x_{k-1})^2\right) \\ p(y_k | x_k) &= \mathcal{N}(y_k | x_k, r) \\ &= \frac{1}{\sqrt{2\pi r}} \exp\left(-\frac{1}{2r}(y_k - x_k)^2\right) \end{aligned} \quad (3.6)$$

which is a discrete-time state space model.

The filtering model (3.1) actually states that the joint prior distribution of the states $(\mathbf{x}_0, \dots, \mathbf{x}_T)$ and the joint likelihood of the measurements $(\mathbf{y}_0, \dots, \mathbf{y}_T)$ are, respectively

$$p(\mathbf{x}_0, \dots, \mathbf{x}_T) = p(\mathbf{x}_0) \prod_{k=1}^T p(\mathbf{x}_k | \mathbf{x}_{k-1}) \quad (3.7)$$

$$p(\mathbf{y}_1, \dots, \mathbf{y}_T | \mathbf{x}_0, \dots, \mathbf{x}_T) = \prod_{k=1}^T p(\mathbf{y}_k | \mathbf{x}_k). \quad (3.8)$$

In principle, for given T we could simply compute the posterior distribution of the states by the Bayes rule:

$$\begin{aligned} p(\mathbf{x}_0, \dots, \mathbf{x}_T | \mathbf{y}_1, \dots, \mathbf{y}_T) &= \frac{p(\mathbf{y}_1, \dots, \mathbf{y}_T | \mathbf{x}_0, \dots, \mathbf{x}_T) p(\mathbf{x}_0, \dots, \mathbf{x}_T)}{p(\mathbf{y}_1, \dots, \mathbf{y}_T)} \\ &\propto p(\mathbf{y}_1, \dots, \mathbf{y}_T | \mathbf{x}_0, \dots, \mathbf{x}_T) p(\mathbf{x}_0, \dots, \mathbf{x}_T). \end{aligned} \quad (3.9)$$

However, this kind of explicit usage of the full Bayes' rule is not feasible in real time applications, because the amount of computations per time step increases when new observations arrive. Thus, this way we could only work with small data sets, because if the amount of data is not bounded (as in real time sensing applications), then at some point of time the computations would become intractable. To cope with real time data we need to have an algorithm which does constant amount of computations per time step.

As discussed in Section 1.2.3, *filtering distributions*, *prediction distributions* and *smoothing distributions* can be computed recursively such that only constant amount of computations is done on each time step. For this reason we shall not consider the full posterior computation at all, but concentrate to the above-mentioned distributions instead. In this chapter, we shall mainly consider computation of the filtering and prediction distributions, and algorithms for computing the smoothing distributions will be considered in the next chapter.

3.1.2 Optimal Filtering Equations

The purpose of *optimal filtering* is to compute the *marginal posterior distribution* of the state \mathbf{x}_k on the time step k given the history of the measurements up to the time step k

$$p(\mathbf{x}_k | \mathbf{y}_{1:k}). \quad (3.10)$$

The fundamental equations of the Bayesian filtering theory are given by the following theorem:

Theorem 3.1 (Bayesian optimal filtering equations). *The recursive equations for computing the predicted distribution $p(\mathbf{x}_k | \mathbf{y}_{1:k-1})$ and the filtering distribution $p(\mathbf{x}_k | \mathbf{y}_{1:k})$ on the time step k are given by the following Bayesian filtering equations:*

- Initialization. *The recursion starts from the prior distribution $p(\mathbf{x}_0)$.*
- Prediction. *The predictive distribution of the state \mathbf{x}_k on time step k given the dynamic model can be computed by the Chapman-Kolmogorov equation*

$$p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) = \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{y}_{1:k-1}) d\mathbf{x}_{k-1}. \quad (3.11)$$

- Update. *Given the measurement \mathbf{y}_k on time step k the posterior distribution of the state \mathbf{x}_k can be computed by the Bayes' rule*

$$p(\mathbf{x}_k | \mathbf{y}_{1:k}) = \frac{1}{Z_k} p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k-1}), \quad (3.12)$$

where the normalization constant Z_k is given as

$$Z_k = \int p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) d\mathbf{x}_k. \quad (3.13)$$

If some of the components of the state are discrete, the corresponding integrals are replaced with summations.

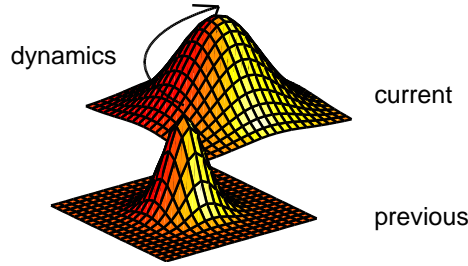


Figure 3.1: Visualization of the prediction step: the prediction propagates the state distribution of the previous measurement step through the dynamic model such that the uncertainties (stochastic terms) in the dynamic model are taken into account.

Proof. The joint distribution of \mathbf{x}_k and \mathbf{x}_{k-1} given $\mathbf{y}_{1:k-1}$ can be computed as

$$\begin{aligned} p(\mathbf{x}_k, \mathbf{x}_{k-1} | \mathbf{y}_{1:k-1}) &= p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{1:k-1}) p(\mathbf{x}_{k-1} | \mathbf{y}_{1:k-1}) \\ &= p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{y}_{1:k-1}), \end{aligned} \quad (3.14)$$

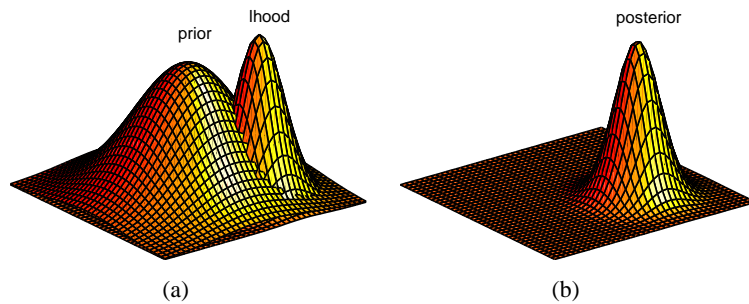


Figure 3.2: Visualization of the update step: (a) Prior distribution from prediction and the likelihood of measurement just before the update step. (b) The posterior distribution after combining the prior and likelihood by Bayes' rule.

where the disappearance of the measurement history $\mathbf{y}_{1:k-1}$ is due to the Markov property of the sequence $\{\mathbf{x}_k, k = 1, 2, \dots\}$. The marginal distribution of \mathbf{x}_k given $\mathbf{y}_{1:k-1}$ can be obtained by integrating the distribution (3.14) over \mathbf{x}_{k-1} , which gives the *Chapman-Kolmogorov equation*

$$p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) = \int p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{y}_{1:k-1}) d\mathbf{x}_{k-1}. \quad (3.15)$$

If \mathbf{x}_{k-1} is discrete, then the above integral is replaced with sum over \mathbf{x}_{k-1} . The distribution of \mathbf{x}_k given \mathbf{y}_k and $\mathbf{y}_{1:k-1}$, that is, given $\mathbf{y}_{1:k}$ can be computed by the *Bayes' rule*

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{y}_{1:k}) &= \frac{1}{Z_k} p(\mathbf{y}_k | \mathbf{x}_k, \mathbf{y}_{1:k-1}) p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) \\ &= \frac{1}{Z_k} p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) \end{aligned} \quad (3.16)$$

where the normalization constant is given by Equation (3.13). The disappearance of the measurement history $\mathbf{y}_{1:k-1}$ in the Equation (3.16) is due to the conditional independence of \mathbf{y}_k from the measurement history, given \mathbf{x}_k . \square

3.1.3 Kalman Filter

The *Kalman filter* (Kalman, 1960b) is the closed form solution to the optimal filtering equations of the discrete-time filtering model, where the dynamic and measurements models are linear Gaussian:

$$\begin{aligned} \mathbf{x}_k &= \mathbf{A}_{k-1} \mathbf{x}_{k-1} + \mathbf{q}_{k-1} \\ \mathbf{y}_k &= \mathbf{H}_k \mathbf{x}_k + \mathbf{r}_k, \end{aligned} \quad (3.17)$$

where $\mathbf{x}_k \in \mathbb{R}^n$ is the state, $\mathbf{y}_k \in \mathbb{R}^m$ is the measurement, $\mathbf{q}_{k-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{k-1})$ is the process noise, $\mathbf{r}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k)$ is the measurement noise and the prior distribution is Gaussian $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{m}_0, \mathbf{P}_0)$. The matrix \mathbf{A}_{k-1} is the transition matrix of the dynamic model and \mathbf{H}_k is the measurement model matrix. In probabilistic terms the model is

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{x}_{k-1}) &= \mathcal{N}(\mathbf{x}_k | \mathbf{A}_{k-1} \mathbf{x}_{k-1}, \mathbf{Q}_{k-1}) \\ p(\mathbf{y}_k | \mathbf{x}_k) &= \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{x}_k, \mathbf{R}_k). \end{aligned} \quad (3.18)$$

Algorithm 3.1 (Kalman filter). *The optimal filtering equations for the linear filtering model (3.17) can be evaluated in closed form and the resulting distributions are Gaussian:*

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) &= \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^-, \mathbf{P}_k^-) \\ p(\mathbf{x}_k | \mathbf{y}_{1:k}) &= \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k) \\ p(\mathbf{y}_k | \mathbf{y}_{1:k-1}) &= \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{m}_k^-, \mathbf{S}_k). \end{aligned} \quad (3.19)$$

The parameters of the distributions above can be computed with the following Kalman filter prediction and update steps:

- The prediction step is

$$\begin{aligned} \mathbf{m}_k^- &= \mathbf{A}_{k-1} \mathbf{m}_{k-1} \\ \mathbf{P}_k^- &= \mathbf{A}_{k-1} \mathbf{P}_{k-1} \mathbf{A}_{k-1}^T + \mathbf{Q}_{k-1}. \end{aligned} \quad (3.20)$$

- The update step is

$$\begin{aligned} \mathbf{v}_k &= \mathbf{y}_k - \mathbf{H}_k \mathbf{m}_k^- \\ \mathbf{S}_k &= \mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k \\ \mathbf{K}_k &= \mathbf{P}_k^- \mathbf{H}_k^T \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k \mathbf{v}_k \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T. \end{aligned} \quad (3.21)$$

The initial state has a given Gaussian prior distribution $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{m}_0, \mathbf{P}_0)$, which also defined the initial mean and covariance.

The Kalman filter equations can be derived as follows:

1. By Lemma A.1 on page 81, the joint distribution of \mathbf{x}_k and \mathbf{x}_{k-1} given $\mathbf{y}_{1:k-1}$ is

$$\begin{aligned} p(\mathbf{x}_{k-1}, \mathbf{x}_k | \mathbf{y}_{1:k-1}) &= p(\mathbf{x}_k | \mathbf{x}_{k-1}) p(\mathbf{x}_{k-1} | \mathbf{y}_{1:k-1}) \\ &= \mathcal{N}(\mathbf{x}_k | \mathbf{A}_{k-1} \mathbf{x}_{k-1}, \mathbf{Q}_{k-1}) \mathcal{N}(\mathbf{x}_{k-1} | \mathbf{m}_{k-1}, \mathbf{P}_{k-1}) \\ &= \mathcal{N} \left(\begin{bmatrix} \mathbf{x}_{k-1} \\ \mathbf{x}_k \end{bmatrix} \middle| \mathbf{m}', \mathbf{P}' \right), \end{aligned} \quad (3.22)$$

where

$$\mathbf{m}' = \begin{pmatrix} \mathbf{m}_{k-1} \\ \mathbf{A}_{k-1} \mathbf{m}_{k-1} \end{pmatrix}, \quad \mathbf{P}' = \begin{pmatrix} \mathbf{P}_{k-1} & \mathbf{P}_{k-1} \mathbf{A}_{k-1}^T \\ \mathbf{A}_{k-1} \mathbf{P}_{k-1} & \mathbf{A}_{k-1} \mathbf{P}_{k-1} \mathbf{A}_{k-1}^T + \mathbf{Q}_{k-1} \end{pmatrix}. \quad (3.23)$$

and the marginal distribution of \mathbf{x}_k is by Lemma A.2

$$p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) = \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^-, \mathbf{P}_k^-), \quad (3.24)$$

where

$$\mathbf{m}_k^- = \mathbf{A}_{k-1} \mathbf{m}_{k-1}, \quad \mathbf{P}_k^- = \mathbf{A}_{k-1} \mathbf{P}_{k-1} \mathbf{A}_{k-1}^T + \mathbf{Q}_{k-1}. \quad (3.25)$$

2. By Lemma A.1, the joint distribution of \mathbf{y}_k and \mathbf{x}_k is

$$\begin{aligned} p(\mathbf{x}_k, \mathbf{y}_k | \mathbf{y}_{1:k-1}) &= p(\mathbf{y}_k | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k-1}) \\ &= \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{x}_k, \mathbf{R}_k) \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^-, \mathbf{P}_k^-) \\ &= \mathcal{N} \left(\begin{bmatrix} \mathbf{x}_k \\ \mathbf{y}_k \end{bmatrix} \middle| \mathbf{m}'' , \mathbf{P}'' \right), \end{aligned} \quad (3.26)$$

where

$$\mathbf{m}'' = \begin{pmatrix} \mathbf{m}_k^- \\ \mathbf{H}_k \mathbf{m}_k^- \end{pmatrix}, \quad \mathbf{P}'' = \begin{pmatrix} \mathbf{P}_k^- & \mathbf{P}_k^- \mathbf{H}_k^T \\ \mathbf{H}_k \mathbf{P}_k^- & \mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k \end{pmatrix}. \quad (3.27)$$

3. By Lemma A.2 the conditional distribution of \mathbf{x}_k is

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{y}_k, \mathbf{y}_{1:k-1}) &= p(\mathbf{x}_k | \mathbf{y}_{1:k}) \\ &= \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k), \end{aligned} \quad (3.28)$$

where

$$\begin{aligned} \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{P}_k^- \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k)^{-1} [\mathbf{y}_k - \mathbf{H}_k \mathbf{m}_k^-] \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{P}_k^- \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^- \mathbf{H}_k^T + \mathbf{R}_k)^{-1} \mathbf{H}_k \mathbf{P}_k^- \end{aligned} \quad (3.29)$$

which can be also written in form (3.21).

The functional form of the Kalman filter equations given here is not the only possible one. In the numerical stability point of view it would be better to work with matrix square roots of covariances instead of plain covariance matrices. The theory and details of implementation of this kind of methods is well covered, for example, in the book of Grewal and Andrews (2001).

Example 3.2 (Kalman filter for Gaussian random walk). *Assume that we are observing measurements y_k of the Gaussian random walk model given in Example 3.1 and we want to estimate the state x_k on each time step. The information obtained up to time step $k - 1$ is summarized by the Gaussian filtering density*

$$p(x_{k-1} | y_{1:k-1}) = \mathcal{N}(x_{k-1} | m_{k-1}, P_{k-1}). \quad (3.30)$$

The Kalman filter prediction and update equations are now given as

$$\begin{aligned}m_k^- &= m_{k-1} \\P_k^- &= P_{k-1} + q \\m_k &= m_k^- + \frac{P_k^-}{P_k^- + r}(y_k - m_k^-) \\P_k &= P_k^- - \frac{(P_k^-)^2}{P_k^- + r}.\end{aligned}\tag{3.31}$$

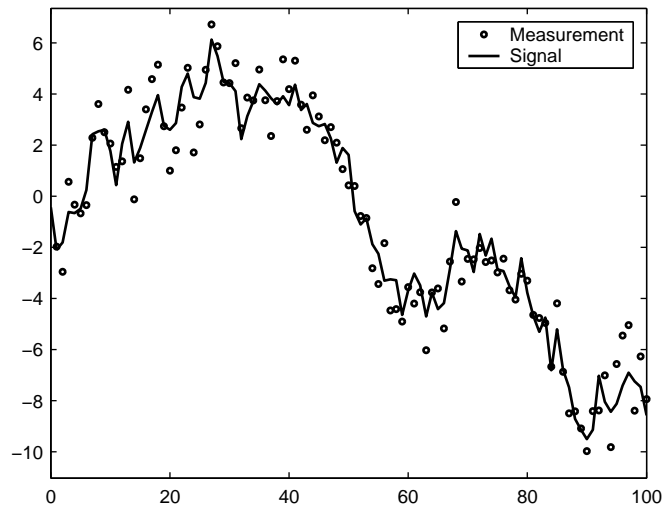


Figure 3.3: Simulated signal and measurements of the Kalman filtering example (Example 3.2).

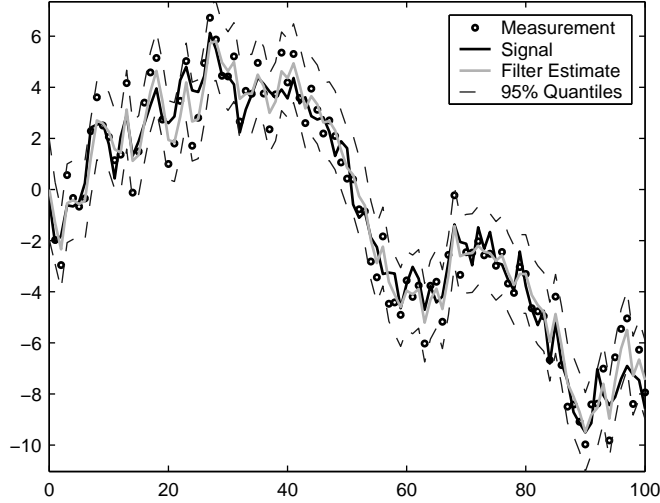


Figure 3.4: Signal, measurements and filtering estimate of the Kalman filtering example (Example 3.2).

3.2 Gaussian Approximation Based Filtering

3.2.1 Linearization of Non-Linear Transforms

Consider the following transformation of a Gaussian random variable \mathbf{x} into another random variable \mathbf{y}

$$\begin{aligned}\mathbf{x} &\sim \mathcal{N}(\mathbf{m}, \mathbf{P}) \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}).\end{aligned}\quad (3.32)$$

where $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$, and $\mathbf{g} : \mathbb{R}^n \mapsto \mathbb{R}^m$ is a general non-linear function. Formally, the probability density of the random variable \mathbf{y} is¹ (see, e.g. Gelman et al., 1995)

$$p(\mathbf{y}) = |\mathbf{J}(\mathbf{y})| \mathcal{N}(\mathbf{g}^{-1}(\mathbf{y}) | \mathbf{m}, \mathbf{P}), \quad (3.33)$$

where $|\mathbf{J}(\mathbf{y})|$ is the determinant of the Jacobian matrix of the inverse transform $\mathbf{g}^{-1}(\mathbf{y})$. However, it is not generally possible to handle this distribution directly, because it is non-Gaussian for all but linear \mathbf{g} .

A first order Taylor series based Gaussian approximation to the distribution of \mathbf{y} can be now formed as follows. If we let $\mathbf{x} = \mathbf{m} + \delta\mathbf{x}$, where $\delta\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{P})$, we can form Taylor series expansion of the function $\mathbf{g}(\cdot)$ as follows:

$$\mathbf{g}(\mathbf{x}) = \mathbf{g}(\mathbf{m} + \delta\mathbf{x}) = \mathbf{g}(\mathbf{m}) + \mathbf{G}_{\mathbf{x}}(\mathbf{m}) \delta\mathbf{x} + \sum_i \frac{1}{2} \delta\mathbf{x}^T \mathbf{G}_{\mathbf{xx}}^{(i)}(\mathbf{m}) \delta\mathbf{x} \mathbf{e}_i + \dots \quad (3.34)$$

¹This actually only applies to invertible $\mathbf{g}(\cdot)$, but it can be easily generalized to the non-invertible case.

where and $\mathbf{G}_x(\mathbf{m})$ is the Jacobian matrix of \mathbf{g} with elements

$$[\mathbf{G}_x(\mathbf{m})]_{j,j'} = \left. \frac{\partial g_j(\mathbf{x})}{\partial x_{j'}} \right|_{\mathbf{x}=\mathbf{m}}. \quad (3.35)$$

and $\mathbf{G}_{xx}^{(i)}(\mathbf{m})$ is the Hessian matrix of $g_i(\cdot)$ evaluated at \mathbf{m} :

$$[\mathbf{G}_{xx}^{(i)}(\mathbf{m})]_{j,j'} = \left. \frac{\partial^2 g_i(\mathbf{x})}{\partial x_j \partial x_{j'}} \right|_{\mathbf{x}=\mathbf{m}}. \quad (3.36)$$

$\mathbf{e}_i = (0 \cdots 0 \ 1 \ 0 \cdots 0)^T$ is a vector with 1 at position i and other elements are zero, that is, it is the unit vector in direction of the coordinate axis i .

The linear approximation can be obtained by approximating the function by the first two terms in the Taylor series:

$$\mathbf{g}(\mathbf{x}) \approx \mathbf{g}(\mathbf{m}) + \mathbf{G}_x(\mathbf{m}) \delta \mathbf{x}. \quad (3.37)$$

Computing the expected value w.r.t. \mathbf{x} gives:

$$\begin{aligned} \mathbf{E}[\mathbf{g}(\mathbf{x})] &\approx \mathbf{E}[\mathbf{g}(\mathbf{m}) + \mathbf{G}_x(\mathbf{m}) \delta \mathbf{x}] \\ &= \mathbf{g}(\mathbf{m}) + \mathbf{G}_x(\mathbf{m}) \mathbf{E}[\delta \mathbf{x}] \\ &= \mathbf{g}(\mathbf{m}). \end{aligned} \quad (3.38)$$

The covariance can be then approximated as

$$\begin{aligned} &\mathbf{E} \left[(\mathbf{g}(\mathbf{x}) - \mathbf{E}[\mathbf{g}(\mathbf{x})]) (\mathbf{g}(\mathbf{x}) - \mathbf{E}[\mathbf{g}(\mathbf{x})])^T \right] \\ &\approx \mathbf{E} \left[(\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{m})) (\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{m}))^T \right] \\ &\approx \mathbf{E} \left[(\mathbf{g}(\mathbf{m}) + \mathbf{G}_x(\mathbf{m}) \delta \mathbf{x} - \mathbf{g}(\mathbf{m})) (\mathbf{g}(\mathbf{m}) + \mathbf{G}_x(\mathbf{m}) \delta \mathbf{x} - \mathbf{g}(\mathbf{m}))^T \right] \\ &= \mathbf{E} \left[(\mathbf{G}_x(\mathbf{m}) \delta \mathbf{x}) (\mathbf{G}_x(\mathbf{m}) \delta \mathbf{x})^T \right] \\ &= \mathbf{G}_x(\mathbf{m}) \mathbf{E} [\delta \mathbf{x} \delta \mathbf{x}^T] \mathbf{G}_x^T(\mathbf{m}) \\ &= \mathbf{G}_x(\mathbf{m}) \mathbf{P} \mathbf{G}_x^T(\mathbf{m}). \end{aligned} \quad (3.39)$$

We are also often interested in the the joint covariance between the variables \mathbf{x} and \mathbf{y} . Approximation to the joint covariance can be achieved by considering the augmented transformation

$$\tilde{\mathbf{g}}(\mathbf{x}) = \begin{pmatrix} \mathbf{x} \\ \mathbf{g}(\mathbf{x}) \end{pmatrix}. \quad (3.40)$$

The resulting mean and covariance are:

$$\begin{aligned} E[\tilde{\mathbf{g}}(\mathbf{x})] &\approx \begin{pmatrix} \mathbf{m} \\ \mathbf{g}(\mathbf{m}) \end{pmatrix} \\ \text{Cov}[\tilde{\mathbf{g}}(\mathbf{x})] &\approx \begin{pmatrix} \mathbf{I} \\ \mathbf{G}_x(\mathbf{m}) \end{pmatrix} \mathbf{P} \begin{pmatrix} \mathbf{I} \\ \mathbf{G}_x(\mathbf{m}) \end{pmatrix}^T \\ &= \begin{pmatrix} \mathbf{P} & \mathbf{P} \mathbf{G}_x^T(\mathbf{m}) \\ \mathbf{G}_x(\mathbf{m}) \mathbf{P} & \mathbf{G}_x(\mathbf{m}) \mathbf{P} \mathbf{G}_x^T(\mathbf{m}) \end{pmatrix}. \end{aligned} \quad (3.41)$$

In the derivation of the extended Kalman filter equations, we need a bit more general transformation of the form

$$\begin{aligned} \mathbf{x} &\sim N(\mathbf{m}, \mathbf{P}) \\ \mathbf{q} &\sim N(\mathbf{0}, \mathbf{Q}) \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}) + \mathbf{q}, \end{aligned} \quad (3.42)$$

where \mathbf{q} is independent of \mathbf{x} . The joint distribution of \mathbf{x} and \mathbf{y} as defined above is now the same as in Equations (3.41) except that the covariance \mathbf{Q} is added to the lower right block of the covariance matrix of $\tilde{\mathbf{g}}(\cdot)$. Thus we get the following algorithm:

Algorithm 3.2 (Linear approximation of additive transform). *The linear approximation based Gaussian approximation to the joint distribution of \mathbf{x} and the transformed random variable $\mathbf{y} = \mathbf{g}(\mathbf{x}) + \mathbf{q}$ where $\mathbf{x} \sim N(\mathbf{m}, \mathbf{P})$ and $\mathbf{q} \sim N(\mathbf{0}, \mathbf{Q})$ is given as*

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{m} \\ \boldsymbol{\mu}_L \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{C}_L \\ \mathbf{C}_L^T & \mathbf{S}_L \end{pmatrix} \right), \quad (3.43)$$

where

$$\begin{aligned} \boldsymbol{\mu}_L &= \mathbf{g}(\mathbf{m}) \\ \mathbf{S}_L &= \mathbf{G}_x(\mathbf{m}) \mathbf{P} \mathbf{G}_x^T(\mathbf{m}) + \mathbf{Q} \\ \mathbf{C}_L &= \mathbf{P} \mathbf{G}_x^T(\mathbf{m}), \end{aligned} \quad (3.44)$$

and $\mathbf{G}_x(\mathbf{m})$ is the Jacobian matrix of \mathbf{g} with respect to \mathbf{x} , evaluated at $\mathbf{x} = \mathbf{m}$ with elements

$$[\mathbf{G}_x(\mathbf{m})]_{j,j'} = \left. \frac{\partial g_j(\mathbf{x})}{\partial x_{j'}} \right|_{\mathbf{x}=\mathbf{m}}. \quad (3.45)$$

Furthermore, in filtering models where the process noise is not additive, we often need to approximate transformations of the form

$$\begin{aligned} \mathbf{x} &\sim N(\mathbf{m}, \mathbf{P}) \\ \mathbf{q} &\sim N(\mathbf{0}, \mathbf{Q}) \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}, \mathbf{q}). \end{aligned} \quad (3.46)$$

where \mathbf{x} and \mathbf{q} are uncorrelated random variables. The mean and covariance can be now computed by substituting the augmented vector (\mathbf{x}, \mathbf{q}) to the vector \mathbf{x} in Equation (3.41). The joint Jacobian matrix can be then written as $\mathbf{G}_{\mathbf{x}, \mathbf{q}} = (\mathbf{G}_{\mathbf{x}} \ \mathbf{G}_{\mathbf{q}})$. Here $\mathbf{G}_{\mathbf{q}}$ is the Jacobian matrix of $\mathbf{g}(\cdot)$ with respect to \mathbf{q} and both the Jacobian matrices are evaluated at $\mathbf{x} = \mathbf{m}, \mathbf{q} = \mathbf{0}$. The approximations to the mean and covariance of the augmented transform as in Equation (3.41) are then given as

$$\begin{aligned} E[\tilde{\mathbf{g}}(\mathbf{x}, \mathbf{q})] &\approx \mathbf{g}(\mathbf{m}, \mathbf{0}) \\ \text{Cov}[\tilde{\mathbf{g}}(\mathbf{x}, \mathbf{q})] &\approx \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{G}_{\mathbf{x}}(\mathbf{m}) & \mathbf{G}_{\mathbf{q}}(\mathbf{m}) \end{pmatrix} \begin{pmatrix} \mathbf{P} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q} \end{pmatrix}^T \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{G}_{\mathbf{x}}(\mathbf{m}) & \mathbf{G}_{\mathbf{q}}(\mathbf{m}) \end{pmatrix}^T \\ &= \begin{pmatrix} \mathbf{P} & \mathbf{P} \mathbf{G}_{\mathbf{x}}^T(\mathbf{m}) \\ \mathbf{G}_{\mathbf{x}}(\mathbf{m}) \mathbf{P} & \mathbf{G}_{\mathbf{x}}(\mathbf{m}) \mathbf{P} \mathbf{G}_{\mathbf{x}}^T(\mathbf{m}) + \mathbf{G}_{\mathbf{q}}(\mathbf{m}) \mathbf{Q} \mathbf{G}_{\mathbf{q}}^T(\mathbf{m}) \end{pmatrix} \end{aligned} \quad (3.47)$$

The approximation above can be formulated as the following algorithm:

Algorithm 3.3 (Linear approximation of non-additive transform). *The linear approximation based Gaussian approximation to the joint distribution of \mathbf{x} and the transformed random variable $\mathbf{y} = \mathbf{g}(\mathbf{x}, \mathbf{q})$ when $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{P})$ and $\mathbf{q} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$ is given as*

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{m} \\ \boldsymbol{\mu}_L \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{C}_L \\ \mathbf{C}_L^T & \mathbf{S}_L \end{pmatrix} \right), \quad (3.48)$$

where

$$\begin{aligned} \boldsymbol{\mu}_L &= \mathbf{g}(\mathbf{m}) \\ \mathbf{S}_L &= \mathbf{G}_{\mathbf{x}}(\mathbf{m}) \mathbf{P} \mathbf{G}_{\mathbf{x}}^T(\mathbf{m}) + \mathbf{G}_{\mathbf{q}}(\mathbf{m}) \mathbf{Q} \mathbf{G}_{\mathbf{q}}^T(\mathbf{m}) \\ \mathbf{C}_L &= \mathbf{P} \mathbf{G}_{\mathbf{x}}^T(\mathbf{m}), \end{aligned} \quad (3.49)$$

and $\mathbf{G}_{\mathbf{x}}(\mathbf{m})$ is the Jacobian matrix of \mathbf{g} with respect to \mathbf{x} , evaluated at $\mathbf{x} = \mathbf{m}, \mathbf{q} = \mathbf{0}$ with elements

$$[\mathbf{G}_{\mathbf{x}}(\mathbf{m})]_{j,j'} = \left. \frac{\partial g_j(\mathbf{x}, \mathbf{q})}{\partial x_{j'}} \right|_{\mathbf{x}=\mathbf{m}, \mathbf{q}=\mathbf{0}}. \quad (3.50)$$

and $\mathbf{G}_{\mathbf{q}}(\mathbf{m})$ is the corresponding Jacobian matrix with respect to \mathbf{q} :

$$[\mathbf{G}_{\mathbf{q}}(\mathbf{m})]_{j,j'} = \left. \frac{\partial g_j(\mathbf{x}, \mathbf{q})}{\partial q_{j'}} \right|_{\mathbf{x}=\mathbf{m}, \mathbf{q}=\mathbf{0}}. \quad (3.51)$$

3.2.2 Extended Kalman Filter

The extended Kalman filter (EKF) (see, e.g., Jazwinski, 1970; Maybeck, 1982a; Bar-Shalom et al., 2001; Grewal and Andrews, 2001) is an extension of the Kalman

filter to non-linear optimal filtering problems. If process and measurement noises can be assumed to be additive, the EKF model can be written as

$$\begin{aligned}\mathbf{x}_k &= \mathbf{f}(\mathbf{x}_{k-1}) + \mathbf{q}_{k-1} \\ \mathbf{y}_k &= \mathbf{h}(\mathbf{x}_k) + \mathbf{r}_k,\end{aligned}\tag{3.52}$$

where $\mathbf{x}_k \in \mathbb{R}^n$ is the state, $\mathbf{y}_k \in \mathbb{R}^m$ is the measurement, $\mathbf{q}_{k-1} \sim \mathcal{N}(0, \mathbf{Q}_{k-1})$ is the Gaussian process noise, $\mathbf{r}_k \sim \mathcal{N}(0, \mathbf{R}_k)$ is the Gaussian measurement noise, $\mathbf{f}(\cdot)$ is the dynamic model function and $\mathbf{h}(\cdot)$ is the measurement model function. The functions \mathbf{f} and \mathbf{h} can also depend on the step number k , but for notational convenience, this dependence has not been explicitly denoted.

The idea of extended Kalman filter is to form Gaussian approximations

$$p(\mathbf{x}_k | \mathbf{y}_{1:k}) \approx \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k),\tag{3.53}$$

to the filtering densities. In EKF this is done by utilizing linear approximations to the non-linearities and the result is

Algorithm 3.4 (Extended Kalman filter I). *The prediction and update steps of the first order additive noise extended Kalman filter are:*

- *Prediction:*

$$\begin{aligned}\mathbf{m}_k^- &= \mathbf{f}(\mathbf{m}_{k-1}) \\ \mathbf{P}_k^- &= \mathbf{F}_x(\mathbf{m}_{k-1}) \mathbf{P}_{k-1} \mathbf{F}_x^T(\mathbf{m}_{k-1}) + \mathbf{Q}_{k-1}.\end{aligned}\tag{3.54}$$

- *Update:*

$$\begin{aligned}\mathbf{v}_k &= \mathbf{y}_k - \mathbf{h}(\mathbf{m}_k^-) \\ \mathbf{S}_k &= \mathbf{H}_x(\mathbf{m}_k^-) \mathbf{P}_k^- \mathbf{H}_x^T(\mathbf{m}_k^-) + \mathbf{R}_k \\ \mathbf{K}_k &= \mathbf{P}_k^- \mathbf{H}_x^T(\mathbf{m}_k^-) \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k \mathbf{v}_k \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.\end{aligned}\tag{3.55}$$

These filtering equations can be derived by repeating the same steps as in derivation of the Kalman filter in Section 3.1.3 and by applying Taylor series approximations on appropriate steps:

1. The joint distribution of \mathbf{x}_k and \mathbf{x}_{k-1} is non-Gaussian, but we can form Gaussian approximation to it by applying the approximation Algorithm 3.2 to the function

$$\mathbf{f}(\mathbf{x}_{k-1}) + \mathbf{q}_{k-1},\tag{3.56}$$

which results in the Gaussian approximation

$$p(\mathbf{x}_{k-1}, \mathbf{x}_k, | \mathbf{y}_{1:k-1}) \approx \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_{k-1} \\ \mathbf{x}_k \end{bmatrix} \middle| \mathbf{m}', \mathbf{P}'\right),\tag{3.57}$$

where

$$\begin{aligned} \mathbf{m}' &= \begin{pmatrix} \mathbf{m}_{k-1} \\ \mathbf{f}(\mathbf{m}_{k-1}) \end{pmatrix} \\ \mathbf{P}' &= \begin{pmatrix} \mathbf{P}_{k-1} & \mathbf{P}_{k-1} \mathbf{F}_x^T \\ \mathbf{F}_x \mathbf{P}_{k-1} & \mathbf{F}_x \mathbf{P}_{k-1} \mathbf{F}_x^T + \mathbf{Q}_{k-1} \end{pmatrix}, \end{aligned} \quad (3.58)$$

and the Jacobian matrix \mathbf{F}_x of $\mathbf{f}(\mathbf{x})$ is evaluated at $\mathbf{x} = \mathbf{m}_{k-1}$. The marginal mean and covariance of \mathbf{x}_k are thus

$$\begin{aligned} \mathbf{m}_k^- &= \mathbf{f}(\mathbf{m}_{k-1}) \\ \mathbf{P}_k^- &= \mathbf{F}_x \mathbf{P}_{k-1} \mathbf{F}_x^T + \mathbf{Q}_{k-1}. \end{aligned} \quad (3.59)$$

2. The joint distribution of \mathbf{y}_k and \mathbf{x}_k is also non-Gaussian, but we can again approximate it by applying the Algorithm 3.2 to the function

$$\mathbf{h}(\mathbf{x}_k) + \mathbf{r}_k. \quad (3.60)$$

We get the approximation

$$p(\mathbf{x}_k, \mathbf{y}_k | \mathbf{y}_{1:k-1}) \approx \mathcal{N} \left(\begin{bmatrix} \mathbf{x}_k \\ \mathbf{y}_k \end{bmatrix} \middle| \mathbf{m}'', \mathbf{P}'' \right), \quad (3.61)$$

where

$$\mathbf{m}'' = \begin{pmatrix} \mathbf{m}_k^- \\ \mathbf{h}(\mathbf{m}_k^-) \end{pmatrix}, \quad \mathbf{P}'' = \begin{pmatrix} \mathbf{P}_k^- & \mathbf{P}_k^- \mathbf{H}_x^T \\ \mathbf{H}_x \mathbf{P}_k^- & \mathbf{H}_x \mathbf{P}_k^- \mathbf{H}_x^T + \mathbf{R}_k \end{pmatrix}, \quad (3.62)$$

and the Jacobian matrix \mathbf{H}_x of $\mathbf{h}(\mathbf{x})$ is evaluated at $\mathbf{x} = \mathbf{m}_k^-$.

3. By Lemma A.2 the conditional distribution of \mathbf{x}_k is approximately

$$p(\mathbf{x}_k | \mathbf{y}_k, \mathbf{y}_{1:k-1}) \approx \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k), \quad (3.63)$$

where

$$\begin{aligned} \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{P}_k^- \mathbf{H}_x^T (\mathbf{H}_x \mathbf{P}_k^- \mathbf{H}_x^T + \mathbf{R}_k)^{-1} [\mathbf{y}_k - \mathbf{h}(\mathbf{m}_k^-)] \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{P}_k^- \mathbf{H}_x^T (\mathbf{H}_x \mathbf{P}_k^- \mathbf{H}_x^T + \mathbf{R}_k)^{-1} \mathbf{H}_x \mathbf{P}_k^- \end{aligned} \quad (3.64)$$

A more general EKF filtering model can be written as

$$\begin{aligned} \mathbf{x}_k &= \mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1}) \\ \mathbf{y}_k &= \mathbf{h}(\mathbf{x}_k, \mathbf{r}_k), \end{aligned} \quad (3.65)$$

where $\mathbf{q}_{k-1} \sim \mathcal{N}(0, \mathbf{Q}_{k-1})$ and $\mathbf{r}_k \sim \mathcal{N}(0, \mathbf{R}_k)$ are the Gaussian process and measurement noises, respectively. Again, the functions \mathbf{f} and \mathbf{h} can also depend on the step number k .

Algorithm 3.5 (Extended Kalman filter II). *The prediction and update steps of the (first order) extended Kalman filter (EKF) are:*

- *Prediction:*

$$\begin{aligned} \mathbf{m}_k^- &= \mathbf{f}(\mathbf{m}_{k-1}, \mathbf{0}) \\ \mathbf{P}_k^- &= \mathbf{F}_x(\mathbf{m}_{k-1}) \mathbf{P}_{k-1} \mathbf{F}_x^T(\mathbf{m}_{k-1}) + \mathbf{F}_q(\mathbf{m}_{k-1}) \mathbf{Q}_{k-1} \mathbf{F}_q^T(\mathbf{m}_{k-1}). \end{aligned} \quad (3.66)$$

- *Update:*

$$\begin{aligned} \mathbf{v}_k &= \mathbf{y}_k - \mathbf{h}(\mathbf{m}_k^-, \mathbf{0}) \\ \mathbf{S}_k &= \mathbf{H}_x(\mathbf{m}_k^-) \mathbf{P}_k^- \mathbf{H}_x^T(\mathbf{m}_k^-) + \mathbf{H}_r(\mathbf{m}_k^-) \mathbf{R}_k \mathbf{H}_r^T(\mathbf{m}_k^-) \\ \mathbf{K}_k &= \mathbf{P}_k^- \mathbf{H}_x^T(\mathbf{m}_k^-) \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k \mathbf{v}_k \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T. \end{aligned} \quad (3.67)$$

where the matrices $\mathbf{F}_x(\mathbf{m})$, $\mathbf{F}_q(\mathbf{m})$, $\mathbf{H}_x(\mathbf{m})$, and $\mathbf{H}_r(\mathbf{m})$, are the Jacobian matrices of \mathbf{f} and \mathbf{h} with respect to state and noise, with elements

$$[\mathbf{F}_x(\mathbf{m})]_{j,j'} = \left. \frac{\partial f_j(\mathbf{x}, \mathbf{q})}{\partial x_{j'}} \right|_{\mathbf{x}=\mathbf{m}, \mathbf{q}=\mathbf{0}} \quad (3.68)$$

$$[\mathbf{F}_q(\mathbf{m})]_{j,j'} = \left. \frac{\partial f_j(\mathbf{x}, \mathbf{q})}{\partial q_{j'}} \right|_{\mathbf{x}=\mathbf{m}, \mathbf{q}=\mathbf{0}} \quad (3.69)$$

$$[\mathbf{H}_x(\mathbf{m})]_{j,j'} = \left. \frac{\partial h_j(\mathbf{x}, \mathbf{r})}{\partial x_{j'}} \right|_{\mathbf{x}=\mathbf{m}, \mathbf{r}=\mathbf{0}} \quad (3.70)$$

$$[\mathbf{H}_r(\mathbf{m})]_{j,j'} = \left. \frac{\partial h_j(\mathbf{x}, \mathbf{r})}{\partial r_{j'}} \right|_{\mathbf{x}=\mathbf{m}, \mathbf{r}=\mathbf{0}}. \quad (3.71)$$

These filtering equations can be derived by repeating the same steps as in the derivation of the extended Kalman filter above, but instead of using the Algorithm 3.2, we use the Algorithm 3.3 for computing the approximations.

In so called second order EKF the non-linearity is approximated by retaining second order terms in Taylor series expansion. The derivation and the resulting equations are straightforward, but due to their complicated appearance, they are not presented here. The equations can be found, for example, in the book of Bar-Shalom et al. (2001).

The advantage of EKF over the other non-linear filtering methods is its relative simplicity compared to its performance. Linearization is very common engineering way of constructing approximations to non-linear systems and thus it is very easy to understand and apply. A disadvantage of it is that because it is based on a

local linear approximation, it will not work in problems with considerable non-linearities. Also the filtering model is restricted in the sense that only Gaussian noise processes are allowed and thus the model cannot contain, for example, discrete valued random variables. The Gaussian restriction also prevents handling of hierarchical models or other models where significantly non-Gaussian distribution models would be needed.

The EKF is also the only filtering algorithm presented in this document, which formally requires the measurement model and dynamic model functions to be differentiable. This as such might be a restriction, but in some cases it might also be simply impossible to compute the required Jacobian matrices, which renders the usage of EKF impossible. And even when the Jacobian matrices exist and could be computed, the actual computation and programming of Jacobian matrices can be quite error prone and hard to debug.

3.2.3 Statistical Linearization of Non-Linear Transforms

In statistically linearized filter (Gelb, 1974) the Taylor series approximation used in the EKF is replaced by statistical linearization. Recall the transformation problem considered in Section 3.2.1, which was stated as

$$\begin{aligned}\mathbf{x} &\sim N(\mathbf{m}, \mathbf{P}) \\ \mathbf{y} &= \mathbf{g}(\mathbf{x}).\end{aligned}$$

In statistical linearization we form a linear approximation to the transformation as follows:

$$\mathbf{g}(\mathbf{x}) \approx \mathbf{b} + \mathbf{A} \delta\mathbf{x}, \quad (3.72)$$

where $\delta\mathbf{x} = \mathbf{x} - \mathbf{m}$, such that the mean squared error is minimized:

$$\text{MSE}(\mathbf{b}, \mathbf{A}) = E[(\mathbf{g}(\mathbf{x}) - \mathbf{b} - \mathbf{A} \delta\mathbf{x})^T (\mathbf{g}(\mathbf{x}) - \mathbf{b} - \mathbf{A} \delta\mathbf{x})]. \quad (3.73)$$

Setting derivatives with respect to \mathbf{b} and \mathbf{A} zero gives

$$\begin{aligned}\mathbf{b} &= E[\mathbf{g}(\mathbf{x})] \\ \mathbf{A} &= E[\mathbf{g}(\mathbf{x}) \delta\mathbf{x}^T] \mathbf{P}^{-1}.\end{aligned} \quad (3.74)$$

In this approximation of transform $\mathbf{g}(\mathbf{x})$, \mathbf{b} is now exactly the mean and the approximate covariance is given as

$$\begin{aligned}E[(\mathbf{g}(\mathbf{x}) - E[\mathbf{g}(\mathbf{x})]) (\mathbf{g}(\mathbf{x}) - E[\mathbf{g}(\mathbf{x})])^T] \\ \approx \mathbf{A} \mathbf{P} \mathbf{A}^T \\ = E[\mathbf{g}(\mathbf{x}) \delta\mathbf{x}^T] \mathbf{P}^{-1} E[\mathbf{g}(\mathbf{x}) \delta\mathbf{x}^T]^T.\end{aligned} \quad (3.75)$$

We may now apply this approximation to the augmented function $\tilde{\mathbf{g}}(\mathbf{x}) = (\mathbf{x}; \mathbf{g}(\mathbf{x}))$ in Equation (3.40) of Section 3.2.1, where we get the approximation

$$\begin{aligned} E[\tilde{\mathbf{g}}(\mathbf{x})] &\approx \begin{pmatrix} \mathbf{m} \\ E[\mathbf{g}(\mathbf{x})] \end{pmatrix} \\ \text{Cov}[\tilde{\mathbf{g}}(\mathbf{x})] &\approx \begin{pmatrix} \mathbf{P} & E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}]^T \\ E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T] & E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T] \mathbf{P}^{-1} E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T]^T \end{pmatrix} \end{aligned} \quad (3.76)$$

We now get the following algorithm corresponding to Algorithm 3.2 in Section 3.2.1:

Algorithm 3.6 (Statistically linearized approximation of additive transform). *The statistical linearization based Gaussian approximation to the joint distribution of \mathbf{x} and the transformed random variable $\mathbf{y} = \mathbf{g}(\mathbf{x}) + \mathbf{q}$ where $\mathbf{x} \sim N(\mathbf{m}, \mathbf{P})$ and $\mathbf{q} \sim N(\mathbf{0}, \mathbf{Q})$ is given as*

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{m} \\ \boldsymbol{\mu}_S \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{C}_S \\ \mathbf{C}_S^T & \mathbf{S}_S \end{pmatrix} \right), \quad (3.77)$$

where

$$\begin{aligned} \boldsymbol{\mu}_S &= E[\mathbf{g}(\mathbf{x})] \\ \mathbf{S}_S &= E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T] \mathbf{P}^{-1} E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T]^T + \mathbf{Q} \\ \mathbf{C}_S &= E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T]^T. \end{aligned} \quad (3.78)$$

The expectations are taken with respect to the distribution of \mathbf{x} .

Applying the same approximation with (\mathbf{x}, \mathbf{q}) in place of \mathbf{x} we obtain the following mean and covariance:

$$\begin{aligned} E[\tilde{\mathbf{g}}(\mathbf{x}, \mathbf{q})] &\approx \begin{pmatrix} \mathbf{m} \\ E[\mathbf{g}(\mathbf{x}, \mathbf{q})] \end{pmatrix} \\ \text{Cov}[\tilde{\mathbf{g}}(\mathbf{x}, \mathbf{q})] &\approx \begin{pmatrix} \mathbf{P} & E[\mathbf{g}(\mathbf{x}, \mathbf{q}) \delta \mathbf{x}^T]^T \\ E[\mathbf{g}(\mathbf{x}, \mathbf{q}) \delta \mathbf{x}^T] & E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T] \mathbf{P}^{-1} E[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T]^T \\ & + E[\mathbf{g}(\mathbf{x}) \mathbf{q}^T] \mathbf{Q}^{-1} E[\mathbf{g}(\mathbf{x}) \mathbf{q}^T]^T \end{pmatrix} \end{aligned} \quad (3.79)$$

Thus we get the following algorithm for non-additive transform as in Algorithm 3.3:

Algorithm 3.7 (Statistically linearized approximation of non-additive transform). *The statistical linearization based Gaussian approximation to the joint distribution of \mathbf{x} and the transformed random variable $\mathbf{y} = \mathbf{g}(\mathbf{x}, \mathbf{q})$ when $\mathbf{x} \sim N(\mathbf{m}, \mathbf{P})$ and $\mathbf{q} \sim N(\mathbf{0}, \mathbf{Q})$ is given as*

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{m} \\ \boldsymbol{\mu}_S \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{C}_S \\ \mathbf{C}_S^T & \mathbf{S}_S \end{pmatrix} \right), \quad (3.80)$$

where

$$\begin{aligned}\boldsymbol{\mu}_S &= \mathbb{E}[\mathbf{g}(\mathbf{x}, \mathbf{q})] \\ \mathbf{S}_S &= \mathbb{E}[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T] \mathbf{P}^{-1} \mathbb{E}[\mathbf{g}(\mathbf{x}) \delta \mathbf{x}^T]^T + \mathbb{E}[\mathbf{g}(\mathbf{x}) \mathbf{q}^T] \mathbf{Q}^{-1} \mathbb{E}[\mathbf{g}(\mathbf{x}) \mathbf{q}^T]^T \\ \mathbf{C}_S &= \mathbb{E}[\mathbf{g}(\mathbf{x}, \mathbf{q}) \delta \mathbf{x}^T]^T.\end{aligned}\quad (3.81)$$

The expectations are taken with respect to variables \mathbf{x} and \mathbf{q} .

3.2.4 Statistically Linearized Filter

Statistically linearized filter (SLF) (Gelb, 1974) or quasi-linear filter (Stengel, 1994) is a Gaussian approximation based filter, which can be applied to the same kind of models as EKF, that is, to models of the form (3.52) or (3.65). The filter is similar to EKF, but the difference is that statistical linearization algorithms 3.6 and 3.7 are used instead of the Taylor series approximations.

Algorithm 3.8 (Statistically linearized filter I). *The prediction and update steps of the additive noise statistically linearized (Kalman) filter are:*

- *Prediction:*

$$\begin{aligned}\mathbf{m}_k^- &= \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1})] \\ \mathbf{P}_k^- &= \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1}) \delta \mathbf{x}_{k-1}^T] \mathbf{P}_{k-1}^{-1} \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1}) \delta \mathbf{x}_{k-1}^T]^T + \mathbf{Q}_{k-1},\end{aligned}\quad (3.82)$$

where $\delta \mathbf{x}_{k-1} = \mathbf{x}_{k-1} - \mathbf{m}_{k-1}$ and the expectations are taken with respect to the variable $\mathbf{x}_{k-1} \sim \mathcal{N}(\mathbf{m}_{k-1}, \mathbf{P}_{k-1})$.

- *Update:*

$$\begin{aligned}\mathbf{v}_k &= \mathbf{y}_k - \mathbb{E}[\mathbf{h}(\mathbf{x}_k)] \\ \mathbf{S}_k &= \mathbb{E}[\mathbf{h}(\mathbf{x}_k) \delta \mathbf{x}_k^T] (\mathbf{P}_k^-)^{-1} \mathbb{E}[\mathbf{h}(\mathbf{x}_k) \delta \mathbf{x}_k^T]^T + \mathbf{R}_k \\ \mathbf{K}_k &= \mathbb{E}[\mathbf{h}(\mathbf{x}_k) \delta \mathbf{x}_k^T]^T \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k \mathbf{v}_k \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.\end{aligned}\quad (3.83)$$

where the expectations are taken with respect to the variable $\mathbf{x}_k \sim \mathcal{N}(\mathbf{m}_k^-, \mathbf{P}_k^-)$.

Algorithm 3.9 (Statistically linearized filter II). *The prediction and update steps of the non-additive statistically linearized (Kalman) filter are:*

- *Prediction:*

$$\begin{aligned}\mathbf{m}_k^- &= \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1})] \\ \mathbf{P}_k^- &= \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1}) \delta \mathbf{x}_{k-1}^T] \mathbf{P}_{k-1}^{-1} \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1}) \delta \mathbf{x}_{k-1}^T]^T \\ &\quad + \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1}) \mathbf{q}_{k-1}^T] \mathbf{Q}_{k-1}^{-1} \mathbb{E}[\mathbf{f}(\mathbf{x}_{k-1}, \mathbf{q}_{k-1}) \mathbf{q}_{k-1}^T]^T,\end{aligned}\quad (3.84)$$

where $\delta \mathbf{x}_{k-1} = \mathbf{x}_{k-1} - \mathbf{m}_{k-1}$ and the expectations are taken with respect to the variables $\mathbf{x}_{k-1} \sim \mathcal{N}(\mathbf{m}_{k-1}, \mathbf{P}_{k-1})$ and $\mathbf{q}_{k-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{k-1})$.

- *Update:*

$$\begin{aligned}
\mathbf{v}_k &= \mathbf{y}_k - \mathbb{E}[\mathbf{h}(\mathbf{x}_k, \mathbf{r}_k)] \\
\mathbf{S}_k &= \mathbb{E}[\mathbf{h}(\mathbf{x}_k, \mathbf{r}_k) \delta \mathbf{x}_k^T] (\mathbf{P}_k^-)^{-1} \mathbb{E}[\mathbf{h}(\mathbf{x}_k, \mathbf{r}_k) \delta \mathbf{x}_k^T]^T \\
&\quad + \mathbb{E}[\mathbf{h}(\mathbf{x}_k, \mathbf{r}_k) \mathbf{r}_k^T] \mathbf{R}_k^{-1} \mathbb{E}[\mathbf{h}(\mathbf{x}_k, \mathbf{r}_k) \mathbf{r}_k^T]^T \\
\mathbf{K}_k &= \mathbb{E}[\mathbf{h}(\mathbf{x}_k, \mathbf{r}_k) \delta \mathbf{x}_k^T]^T \mathbf{S}_k^{-1} \\
\mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k \mathbf{v}_k \\
\mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.
\end{aligned} \tag{3.85}$$

where the expectations are taken with respect to variables $\mathbf{x}_k \sim \mathcal{N}(\mathbf{m}_k^-, \mathbf{P}_k^-)$ and $\mathbf{r}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k)$.

Both the filters above can be derived by following the derivation of the EKF in Section 3.2.2 and by utilizing the statistical linearization approximations instead of linear approximations on appropriate steps.

The advantage of SLF over EKF is that it is more global approximation than EKF, because the linearization is not only based on the local region around the mean but on a whole range of function values. The non-linearities also do not have to be differentiable nor do we need to derive their Jacobian matrices. The clear disadvantage is that the certain expected values of the non-linear functions have to be computed in closed form. Naturally, it is not possible for all functions. Fortunately, the expected values involved are of such type that one is likely to find many of them tabulated in older physics and control engineering books.

3.2.5 Unscented Transform

The *unscented transform* (UT) (Julier and Uhlmann, 1995; Julier et al., 2000) is a relatively recent numerical method, which can be also used for approximating the joint distribution of random variables \mathbf{x} and \mathbf{y} defined as

$$\begin{aligned}
\mathbf{x} &\sim \mathcal{N}(\mathbf{m}, \mathbf{P}) \\
\mathbf{y} &= \mathbf{g}(\mathbf{x}).
\end{aligned}$$

However, the philosophy in UT differs from the linearization and statistical linearization in the sense that it tries to directly approximate the mean and covariance of the target distribution instead of trying to approximate the non-linear function (Julier and Uhlmann, 1995).

The idea of UT is to form a fixed number of deterministically chosen sigma-points, which capture the mean and covariance of the original distribution of \mathbf{x} exactly. These sigma-points are then propagated through the non-linearity and the mean and covariance of the transformed variable are estimated from them. Note that although the unscented transform resembles Monte Carlo estimation the approaches are significantly different, because in UT the sigma points are selected

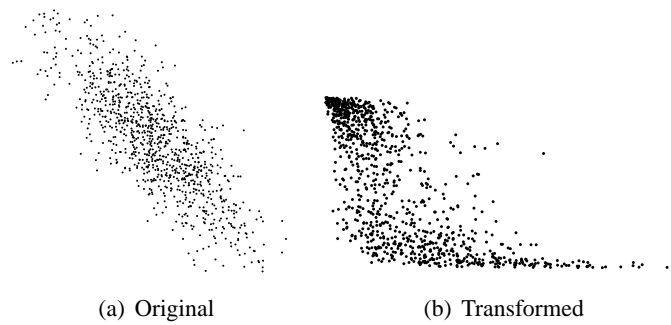


Figure 3.5: Example of applying a non-linear transformation to a random variable on the left, which results in the random variable on the right.

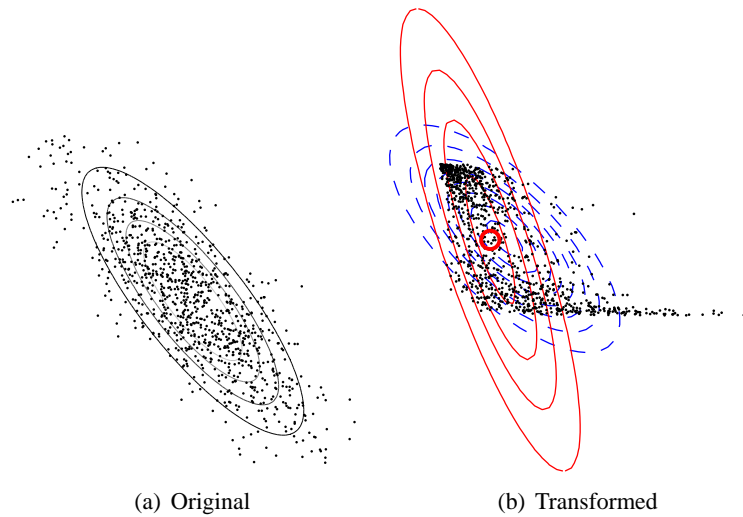


Figure 3.6: Illustration of linearization based (EKF) approximation to the transformation in Figure 3.5. The Gaussian approximation is formed by calculating the curvature at the mean, which results in bad approximation further from the mean. The true distribution is presented by the blue dotted line and the red solid line is the approximation.

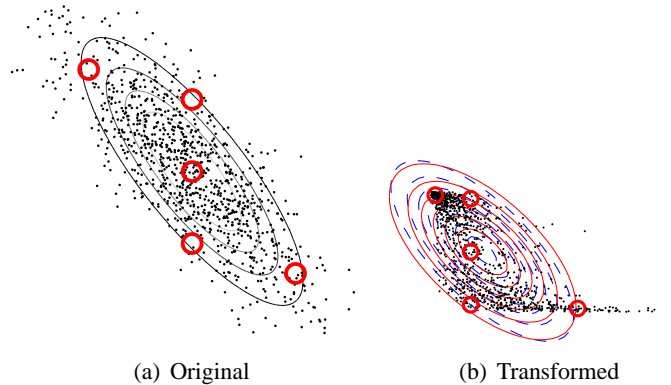


Figure 3.7: Illustration of unscented transform based (UKF) approximation to the transformation in Figure 3.5. The Gaussian approximation is formed by propagating the sigma points through the non-linearity and the mean and covariance are estimated from the transformed sigma points. The true distribution is presented by the blue dotted line and the red solid line is the approximation.

deterministically (Julier and Uhlmann, 2004). The difference between linear approximation and UT is illustrated in Figures 3.5, 3.6 and 3.7.

The *unscented transform* forms the Gaussian approximation with the following procedure:

1. Form the matrix of sigma points \mathbf{X} as

$$\mathbf{X} = [\mathbf{m} \quad \cdots \quad \mathbf{m}] + \sqrt{n + \lambda} [\mathbf{0} \quad \sqrt{\mathbf{P}} \quad -\sqrt{\mathbf{P}}],$$

where λ is a scaling parameter, which is defined in terms of algorithm parameters α and κ as follows:

$$\lambda = \alpha^2 (n + \kappa) - n. \quad (3.86)$$

The parameters α and κ determine the spread of the sigma points around the mean (Wan and Van der Merwe, 2001). The matrix square root denotes a matrix such that $\sqrt{\mathbf{P}} \sqrt{\mathbf{P}}^T = \mathbf{P}$. The sigma points are the columns of the sigma point matrix.

2. Propagate the sigma points through the non-linear function $\mathbf{g}(\cdot)$:

$$\mathbf{Y}_i = \mathbf{g}(\mathbf{X}_i), \quad i = 1 \dots 2n + 1,$$

where \mathbf{X}_i and \mathbf{Y}_i denote the i th columns of matrices \mathbf{X} and \mathbf{Y} , respectively.

3. Estimates of the mean and covariance of the transformed variable can be

computed from the sigma points as follows:

$$\begin{aligned} E[\mathbf{g}(\mathbf{x})] &\approx \sum_i W_{i-1}^{(m)} \mathbf{Y}_i \\ \text{Cov}[\mathbf{g}(\mathbf{x})] &\approx \sum_i W_{i-1}^{(c)} (\mathbf{Y}_i - \boldsymbol{\mu})(\mathbf{Y}_i - \boldsymbol{\mu})^T \end{aligned} \quad (3.87)$$

where the constant weights $W_i^{(m)}$ and $W_i^{(c)}$ are given as follows (Wan and Van der Merwe, 2001):

$$\begin{aligned} W_0^{(m)} &= \lambda/(n + \lambda) \\ W_0^{(c)} &= \lambda/(n + \lambda) + (1 - \alpha^2 + \beta) \\ W_i^{(m)} &= 1/\{2(n + \lambda)\}, \quad i = 1, \dots, 2n \\ W_i^{(c)} &= 1/\{2(n + \lambda)\}, \quad i = 1, \dots, 2n, \end{aligned} \quad (3.88)$$

and β is an additional algorithm parameter, which can be used for incorporating prior information on the (non-Gaussian) distribution of \mathbf{x} (Wan and Van der Merwe, 2001). Note that the indexing starts from zero, because originally the sigma points were numbered starting from zero instead of starting from one as we do here.

If we apply the unscented transform to the augmented function $\tilde{\mathbf{g}}(\mathbf{x}) = (\mathbf{x}, \mathbf{g}(\mathbf{x}))$, we simply get the set of sigma points, where the sigma points \mathbf{X}_i and \mathbf{Y}_i have been concatenated to the same vectors. Thus, also forming approximation to joint distribution \mathbf{x} and $\mathbf{g}(\mathbf{x}) + \mathbf{q}$ is straightforward and the result is:

Algorithm 3.10 (Unscented approximation of additive transform). *The unscented transform approximation based Gaussian approximation to the joint distribution of \mathbf{x} and the transformed random variable $\mathbf{y} = \mathbf{g}(\mathbf{x}) + \mathbf{q}$ where $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{P})$ and $\mathbf{q} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$ is given as*

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{m} \\ \boldsymbol{\mu}_U \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{C}_U \\ \mathbf{C}_U^T & \mathbf{S}_U \end{pmatrix} \right), \quad (3.89)$$

where the submatrices can be computed as follows:

1. Form the matrix of sigma points \mathbf{X} as

$$\mathbf{X} = [\mathbf{m} \quad \cdots \quad \mathbf{m}] + \sqrt{n + \lambda} [0 \quad \sqrt{\mathbf{P}} \quad -\sqrt{\mathbf{P}}],$$

where the parameters are as defined above.

2. Propagate the sigma points through the non-linear function $\mathbf{g}(\cdot)$:

$$\mathbf{Y}_i = \mathbf{g}(\mathbf{X}_i), \quad i = 1 \dots 2n + 1,$$

where \mathbf{X}_i and \mathbf{Y}_i denote the i th columns of matrices \mathbf{X} and \mathbf{Y} , respectively.

3. The submatrices are then given as:

$$\begin{aligned}\boldsymbol{\mu}_U &= \sum_i W_{i-1}^{(m)} \mathbf{Y}_i \\ \mathbf{S}_U &= \sum_i W_{i-1}^{(c)} (\mathbf{Y}_i - \boldsymbol{\mu}_U) (\mathbf{Y}_i - \boldsymbol{\mu}_U)^T + \mathbf{Q} \\ \mathbf{C}_U &= \sum_i W_{i-1}^{(c)} (\mathbf{X}_i - \mathbf{m}) (\mathbf{Y}_i - \boldsymbol{\mu}_U)^T,\end{aligned}\quad (3.90)$$

where the constant weights $W_i^{(m)}$ and $W_i^{(c)}$ were defined above.

The unscented transform approximation to a transformation of the form $\mathbf{y} = \mathbf{g}(\mathbf{x}, \mathbf{q})$ can be derived by considering the augmented random variable $\tilde{\mathbf{x}} = (\mathbf{x}, \mathbf{q})$ as the input variable. The resulting algorithm is:

Algorithm 3.11 (Unscented approximation of non-additive transform). *The unscented transform based Gaussian approximation to the joint distribution of \mathbf{x} and the transformed random variable $\mathbf{y} = \mathbf{g}(\mathbf{x}, \mathbf{q})$ when $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \mathbf{P})$ and $\mathbf{q} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$ is given as*

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{m} \\ \boldsymbol{\mu}_U \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{C}_U \\ \mathbf{C}_U^T & \mathbf{S}_U \end{pmatrix} \right), \quad (3.91)$$

where the submatrices can be computed as follows. Let the dimensionalities of \mathbf{x} and \mathbf{q} be n_x and n_q , respectively, and let $n = n_x + n_q$.

1. Form the matrix of sigma points of the augmented random variable $\tilde{\mathbf{x}} = (\mathbf{x}, \mathbf{q})$

$$\tilde{\mathbf{X}} = (\tilde{\mathbf{m}} \quad \cdots \quad \tilde{\mathbf{m}}) + \sqrt{n + \lambda} \begin{pmatrix} \mathbf{0} & \sqrt{\tilde{\mathbf{P}}} & -\sqrt{\tilde{\mathbf{P}}} \end{pmatrix}.$$

where

$$\tilde{\mathbf{m}} = \begin{pmatrix} \mathbf{m} \\ \mathbf{0} \end{pmatrix} \quad \tilde{\mathbf{P}} = \begin{pmatrix} \mathbf{P} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q} \end{pmatrix}.$$

2. Propagate the sigma points through the function:

$$\tilde{\mathbf{Y}}_i = \mathbf{g}(\tilde{\mathbf{X}}_i^x, \tilde{\mathbf{X}}_i^q), \quad i = 1 \dots 2n + 1,$$

where $\tilde{\mathbf{X}}_i^x$ and $\tilde{\mathbf{X}}_i^q$ denote the parts of the augmented sigma point i , which correspond to \mathbf{x} and \mathbf{q} , respectively.

3. Compute the predicted mean $\boldsymbol{\mu}_U$, the predicted covariance \mathbf{S}_U and the cross-covariance \mathbf{C}_U :

$$\begin{aligned}\boldsymbol{\mu}_U &= \sum_i W_{i-1}^{(m)} \tilde{\mathbf{Y}}_i \\ \mathbf{S}_U &= \sum_i W_{i-1}^{(c)} (\tilde{\mathbf{Y}}_i - \boldsymbol{\mu}_U) (\tilde{\mathbf{Y}}_i - \boldsymbol{\mu}_U)^T \\ \mathbf{C}_U &= \sum_i W_{i-1}^{(c)} (\tilde{\mathbf{X}}_i^x - \mathbf{m}) (\tilde{\mathbf{Y}}_i - \boldsymbol{\mu}_U)^T,\end{aligned}$$

where the definitions of the weights $W_i^{(m)}$ and $W_i^{(c)}$ are as above.

3.2.6 Unscented Kalman Filter

The *unscented Kalman filter* (UKF) (Julier et al., 1995; Julier and Uhlmann, 2004; Wan and Van der Merwe, 2001) is an optimal filtering algorithm that utilizes the unscented transform and can be used for approximating the filtering distribution of models having the same form as with EKF and SLF, that is, models of the form (3.52) or (3.65). As EKF and SLF, UKF forms a Gaussian approximation to the filtering distribution:

$$p(\mathbf{x}_k \mid \mathbf{y}_1, \dots, \mathbf{y}_k) \approx \mathcal{N}(\mathbf{x}_k \mid \mathbf{m}_k, \mathbf{P}_k), \quad (3.92)$$

where \mathbf{m}_k and \mathbf{P}_k are the mean and covariance computed by the algorithm.

Algorithm 3.12 (unscented Kalman filter). *In the additive form unscented Kalman filter (UKF) algorithm, which can be applied to additive models of the form (3.52), the following operations are performed on each measurement step $k = 1, 2, 3, \dots$:*

1. Prediction step:

(a) Form the matrix of sigma points:

$$\mathbf{X}_{k-1} = [\mathbf{m}_{k-1} \quad \cdots \quad \mathbf{m}_{k-1}] + \sqrt{n + \lambda} \begin{bmatrix} 0 & \sqrt{\mathbf{P}_{k-1}} & -\sqrt{\mathbf{P}_{k-1}} \end{bmatrix}. \quad (3.93)$$

(b) Propagate the sigma points through the dynamic model:

$$\hat{\mathbf{X}}_{k,i} = \mathbf{f}(\mathbf{X}_{k-1,i}), \quad i = 1 \dots 2n + 1. \quad (3.94)$$

(c) Compute the predicted mean \mathbf{m}_k^- and the predicted covariance \mathbf{P}_k^- :

$$\begin{aligned} \mathbf{m}_k^- &= \sum_i W_{i-1}^{(m)} \hat{\mathbf{X}}_{k,i} \\ \mathbf{P}_k^- &= \sum_i W_{i-1}^{(c)} (\hat{\mathbf{X}}_{k,i} - \mathbf{m}_k^-) (\hat{\mathbf{X}}_{k,i} - \mathbf{m}_k^-)^T + \mathbf{Q}_{k-1}. \end{aligned} \quad (3.95)$$

2. Update step:

(a) Form the matrix of sigma points:

$$\mathbf{X}_k^- = [\mathbf{m}_k^- \quad \cdots \quad \mathbf{m}_k^-] + \sqrt{n + \lambda} \begin{bmatrix} 0 & \sqrt{\mathbf{P}_k^-} & -\sqrt{\mathbf{P}_k^-} \end{bmatrix}. \quad (3.96)$$

(b) Propagate sigma points through the measurement model:

$$\hat{\mathbf{Y}}_{k,i} = \mathbf{h}(\mathbf{X}_{k,i}^-), \quad i = 1 \dots 2n + 1. \quad (3.97)$$

- (c) Compute the predicted mean $\boldsymbol{\mu}_k$, the predicted covariance of the measurement \mathbf{S}_k , and the cross-covariance of the state and measurement \mathbf{C}_k :

$$\begin{aligned}\boldsymbol{\mu}_k &= \sum_i W_{i-1}^{(m)} \hat{\mathbf{Y}}_{k,i} \\ \mathbf{S}_k &= \sum_i W_{i-1}^{(c)} (\hat{\mathbf{Y}}_{k,i} - \boldsymbol{\mu}_k) (\hat{\mathbf{Y}}_{k,i} - \boldsymbol{\mu}_k)^T + \mathbf{R}_k \\ \mathbf{C}_k &= \sum_i W_{i-1}^{(c)} (\mathbf{X}_{k,i}^- - \mathbf{m}_k^-) (\hat{\mathbf{Y}}_{k,i} - \boldsymbol{\mu}_k)^T.\end{aligned}\quad (3.98)$$

- (d) Compute the filter gain \mathbf{K}_k and the filtered state mean \mathbf{m}_k and covariance \mathbf{P}_k , conditional to the measurement \mathbf{y}_k :

$$\begin{aligned}\mathbf{K}_k &= \mathbf{C}_k \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k [\mathbf{y}_k - \boldsymbol{\mu}_k] \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T.\end{aligned}\quad (3.99)$$

The filtering equations above can be derived in analogous manner to EKF equations, but the unscented transform based approximations are used instead of the linear approximations.

The non-additive form of UKF (Julier and Uhlmann, 2004) can be derived by augmenting the process or measurement noises with the state vector and applying UT approximation to that. Alternatively, one can first augment the state vector with process noise, then approximate the prediction step and after that do the same with measurement noise on the update step. The different algorithms and ways of doing this in practice are analyzed in article (Wu et al., 2005). Because of the various alternative forms and complicated appearance of each of these, the reader is encouraged to check the augmented form filtering equations from articles (Julier and Uhlmann, 2004; Wu et al., 2005) and references therein.

The advantage of the UKF over EKF is that UKF is not based on local linear approximation, but uses a bit further points in approximating the non-linearity. As discussed in Julier and Uhlmann (2004) the unscented transform is able to capture the higher order moments caused by the non-linear transform better than the Taylor series based approximations. The dynamic and model functions are also not required to be formally differentiable nor their Jacobian matrices need to be computed. The advantage of UKF over SLF is that in UKF there is no need to compute any expected values in closed form, only evaluations of the dynamic and measurement models are needed. However, the accuracy of UKF cannot be expected to be as good as of SLF, because SLF try uses larger area in the approximation, whereas UKF only selects fixed number of points on the area. The disadvantage over EKF is that UKF often requires slightly more computational operations than EKF.

The UKF can be interpreted to belong to a wider class of filters called sigma-point filters (van der Merwe and Wan, 2003), which also includes other types of

filters such as central differences Kalman filter (CDKF), Gauss-Hermite Kalman filter (GHKF) and a few others (Ito and Xiong, 2000; Wu et al., 2006; Nørgaard et al., 2000; Arasaratnam and Haykin, 2009). The classification to sigma-point methods by van der Merwe and Wan (2003) is based on interpreting the methods as special cases of (weighted) statistical linear regression (Lefebvre et al., 2002).

As discussed in (van der Merwe and Wan, 2003), statistical linearization is closely related to sigma-point approximations, because they both are related to statistical linear regression. However, it is important to note that the statistical linear regression (Lefebvre et al., 2002) which is the basis of sigma-point framework (van der Merwe and Wan, 2003) is not exactly equivalent to statistical linearization (Gelb, 1974) as sometimes is claimed. The statistical linear regression can be considered as a discrete approximation to statistical linearization.

3.2.7 Gaussian Moment Matching

One way to unify the Taylor series, statistical linearization and unscented transform based approaches is to think all of them as approximations to the moment integrals:

$$\begin{aligned}\boldsymbol{\mu}_M &= \int \mathbf{g}(\mathbf{x}) N(\mathbf{x} | \mathbf{m}, \mathbf{P}) d\mathbf{x} \\ \mathbf{S}_M &= \int (\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_M) (\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_M)^T N(\mathbf{x} | \mathbf{m}, \mathbf{P}) d\mathbf{x} \\ \mathbf{C}_M &= \int (\mathbf{x} - \mathbf{m}) (\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_M)^T N(\mathbf{x} | \mathbf{m}, \mathbf{P}) d\mathbf{x}.\end{aligned}$$

If we can compute these, a straight-forward way to form the Gaussian approximation for (\mathbf{x}, \mathbf{y}) is to simply match the moments of the distributions, which gives the following algorithm:

Algorithm 3.13 (Gaussian moment matching of additive transform). *The moment matching based Gaussian approximation to the joint distribution of \mathbf{x} and the transformed random variable $\mathbf{y} = \mathbf{g}(\mathbf{x}) + \mathbf{q}$ where $\mathbf{x} \sim N(\mathbf{m}, \mathbf{P})$ and $\mathbf{q} \sim N(\mathbf{0}, \mathbf{Q})$ is given as*

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N \left(\begin{pmatrix} \mathbf{m} \\ \boldsymbol{\mu}_M \end{pmatrix}, \begin{pmatrix} \mathbf{P} & \mathbf{C}_M \\ \mathbf{C}_M^T & \mathbf{S}_M \end{pmatrix} \right), \quad (3.100)$$

where

$$\begin{aligned}\boldsymbol{\mu}_M &= \int \mathbf{g}(\mathbf{x}) N(\mathbf{x} | \mathbf{m}, \mathbf{P}) d\mathbf{x} \\ \mathbf{S}_M &= \int (\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_M) (\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_M)^T N(\mathbf{x} | \mathbf{m}, \mathbf{P}) d\mathbf{x} + \mathbf{Q} \\ \mathbf{C}_M &= \int (\mathbf{x} - \mathbf{m}) (\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_M)^T N(\mathbf{x} | \mathbf{m}, \mathbf{P}) d\mathbf{x}.\end{aligned} \quad (3.101)$$

The non-additive case can be handled in analogous manner. It is now easy to check by substituting the approximation $\mathbf{g}(\mathbf{x}) = \mathbf{g}(\mathbf{m}) + \mathbf{G}_x(\mathbf{m})(\mathbf{x} - \mathbf{m})$

to the above expression that in the linear case the integrals indeed reduce to the linear approximations in the Algorithm 3.2. And the same applies to statistical linearization. However, many other approximations can also be interpreted as such approximations as is discussed in the next section.

3.2.8 Gaussian Assumed Density Filter

If we replace the linear approximations in EKF with the moment matching approximations in the previous section, we get the following *Gaussian assumed density filter* (ADF) which is also called *Gaussian filter* (Maybeck, 1982a; Ito and Xiong, 2000; Wu et al., 2006):

Algorithm 3.14 (Gaussian assumed density filter). *The prediction and update steps of the additive noise Gaussian assumed density (Kalman) filter are:*

- *Prediction:*

$$\begin{aligned} \mathbf{m}_k^- &= \int \mathbf{f}(\mathbf{x}_{k-1}) \mathcal{N}(\mathbf{x}_{k-1} | \mathbf{m}_{k-1}, \mathbf{P}_{k-1}) d\mathbf{x}_{k-1} \\ \mathbf{P}_k^- &= \int (\mathbf{f}(\mathbf{x}_{k-1}) - \mathbf{m}_k^-) (\mathbf{f}(\mathbf{x}_{k-1}) - \mathbf{m}_k^-)^T \\ &\quad \times \mathcal{N}(\mathbf{x}_{k-1} | \mathbf{m}_{k-1}, \mathbf{P}_{k-1}) d\mathbf{x}_{k-1} + \mathbf{Q}_{k-1}. \end{aligned} \quad (3.102)$$

- *Update:*

$$\begin{aligned} \boldsymbol{\mu}_k &= \int \mathbf{h}(\mathbf{x}_k) \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^-, \mathbf{P}_k^-) d\mathbf{x}_k \\ \mathbf{S}_k &= \int (\mathbf{h}(\mathbf{x}_k) - \boldsymbol{\mu}_k) (\mathbf{h}(\mathbf{x}_k) - \boldsymbol{\mu}_k)^T \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^-, \mathbf{P}_k^-) d\mathbf{x}_k + \mathbf{R}_k \\ \mathbf{C}_k &= \int (\mathbf{x}_k - \mathbf{m}_k^-) (\mathbf{h}(\mathbf{x}_k) - \boldsymbol{\mu}_k)^T \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^-, \mathbf{P}_k^-) d\mathbf{x}_k \\ \mathbf{K}_k &= \mathbf{C}_k \mathbf{S}_k^{-1} \\ \mathbf{m}_k &= \mathbf{m}_k^- + \mathbf{K}_k (\mathbf{y}_k - \boldsymbol{\mu}_k) \\ \mathbf{P}_k &= \mathbf{P}_k^- - \mathbf{K}_k \mathbf{S}_k \mathbf{K}_k^T. \end{aligned} \quad (3.103)$$

The advantage of the moment matching formulation is that it enables usage of many well known numerical integration methods such as Gauss-Hermite quadratures, cubature rules and central difference based methods (Ito and Xiong, 2000; Wu et al., 2006; Nørgaard et al., 2000; Arasaratnam and Haykin, 2009). The unscented transformation can also be interpreted as an approximation to these integrals (Wu et al., 2006).

One interesting way to approximate the integrals is to use the Bayes-Hermite quadrature (O'Hagan, 1991), which is based of fitting a Gaussian process regression model to the non-linear functions on finite set of training points. This approach

is used in the Gaussian process filter of Deisenroth et al. (2009). It is also possible to approximate the integrals by Monte Carlo integration, which is the approach used in Monte Carlo Kalman Filter (MCKF) of Kotecha and Djuric (2003).

3.3 Monte Carlo Approximations

3.3.1 Principles and Motivation of Monte Carlo

Within statistical methods in engineering and science, as well as in optimal filtering, it is often necessary to evaluate expectations in form

$$E[\mathbf{g}(\mathbf{x})] = \int \mathbf{g}(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}, \quad (3.104)$$

where $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ in an arbitrary function and $p(\mathbf{x})$ is the probability density of \mathbf{x} . Now the problem is that such an integral can be evaluated in closed form only in a few special cases and generally, numerical methods have to be used.

Monte Carlo methods provide a numerical method for calculating integrals of the form (3.104). Monte Carlo refers to general class of methods, where closed form computation of statistical quantities is replaced by drawing samples from the distribution and estimating the quantities by sample averages.

In (perfect) Monte Carlo approximation, we draw independent random samples from $\mathbf{x}^{(i)} \sim p(\mathbf{x})$ and estimate the expectation as

$$E[\mathbf{g}(\mathbf{x})] \approx \frac{1}{N} \sum_i \mathbf{g}(\mathbf{x}^{(i)}). \quad (3.105)$$

Thus Monte Carlo methods approximate the target density by a set of samples that are distributed according to the target density. Figure 3.8 represents a two dimensional Gaussian distribution and its Monte Carlo representation.

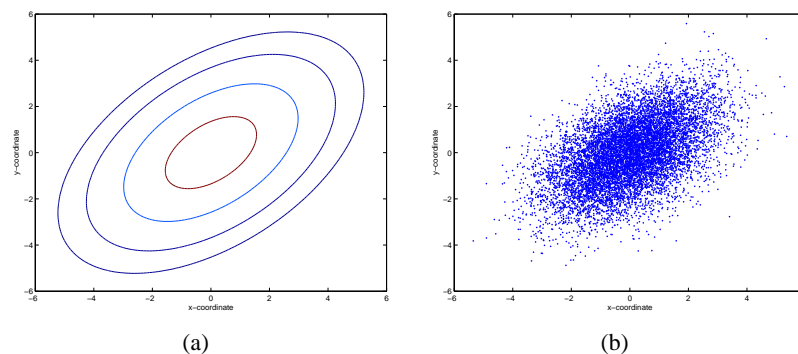


Figure 3.8: (a) Two dimensional Gaussian density. (b) Monte Carlo representation of the same Gaussian density.

The convergence of Monte Carlo approximation is guaranteed by Central Limit Theorem (CLT) (see, e.g., Liu, 2001) and the error term is $O(N^{-1/2})$, regardless of dimensionality of \mathbf{x} . This invariance of dimensionality is unique to Monte Carlo methods and makes them superior to practically all other numerical methods when dimensionality of \mathbf{x} is considerable. At least in theory, not necessarily in practice.

In Bayesian inference the target distribution is typically the posterior distribution $p(\mathbf{x} | \mathbf{y}_1, \dots, \mathbf{y}_n)$ and it is assumed that it is easier to draw (weighted) samples from the distribution than to compute, for example, integrals of the form (3.104). This, indeed, often happens to be the case.

3.3.2 Importance Sampling

It is not always possible to obtain samples directly from $p(\mathbf{x})$ due to its complicated formal appearance. In *importance sampling* (IS) (see, e.g., Liu, 2001) we use approximate distribution called importance distribution $\pi(\mathbf{x})$, which we can easily draw samples from. Having samples $\mathbf{x}^{(i)} \sim \pi(\mathbf{x})$ we can approximate the expectation integral (3.104) as

$$E[\mathbf{g}(\mathbf{x})] \approx \frac{1}{N} \sum_i \frac{\mathbf{g}(\mathbf{x}^{(i)}) p(\mathbf{x}^{(i)})}{\pi(\mathbf{x}^{(i)})}. \quad (3.106)$$

Figure 3.9 illustrates the idea of importance sampling. We sample from the importance distribution, which is an approximation to the target distribution. Because the distribution of samples is not exact, we need to correct the approximation by associating a weight to each of the samples.

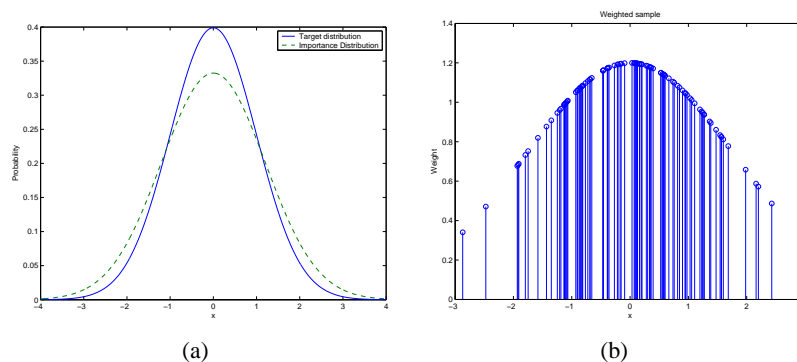


Figure 3.9: (a) Importance distribution approximates the target distribution (b) Weights are associated to each of the samples to correct the approximation.

The disadvantage of this direct importance sampling is that we should be able to evaluate $p(\mathbf{x}^{(i)})$ in order to use it directly. But the problem is that we often do not know the normalization constant of $p(\mathbf{x}^{(i)})$, because evaluation of it would require

evaluation of an integral with comparable complexity to the expectation integral itself. In importance sampling we often use an approximation, where we define unnormalized weights as

$$w_i = \frac{p(\mathbf{x}^{(i)})}{\pi(\mathbf{x}^{(i)})}. \quad (3.107)$$

and approximate the expectation as

$$E[\mathbf{g}(\mathbf{x})] \approx \frac{\sum_i \mathbf{g}(\mathbf{x}^{(i)}) w_i}{\sum_i w_i}, \quad (3.108)$$

which has the fortunate property that we do not have to know the normalization constant of $p(\mathbf{x})$.

3.4 Particle Filtering

3.4.1 Sequential Importance Sampling

Sequential importance sampling (SIS) (see, e.g., Doucet et al., 2001) is a sequential version of importance sampling. It is based on fact that we can evaluate the importance distribution for states \mathbf{x}_k on each time step k recursively as follows:

$$\pi(\mathbf{x}_{0:k} | \mathbf{y}_{1:k}) = \pi(\mathbf{x}_k | \mathbf{x}_{0:k-1}, \mathbf{y}_{1:k}) \pi(\mathbf{x}_{0:k-1} | \mathbf{y}_{1:k-1}) \quad (3.109)$$

Thus, we can also evaluate the (unnormalized) importance weights recursively:

$$\tilde{w}_k^{(i)} \propto \tilde{w}_{k-1}^{(i)} \frac{p(\mathbf{y}_k | \mathbf{x}_k^{(i)}) p(\mathbf{x}_k^{(i)} | \mathbf{x}_{k-1}^{(i)})}{\pi(\mathbf{x}_k^{(i)} | \mathbf{x}_{0:k-1}^{(i)}, \mathbf{y}_{1:k})} \quad (3.110)$$

The SIS algorithm can be used for generating Monte Carlo approximations to filtering distributions of generic state space models of the form

$$\begin{aligned} \mathbf{x}_k &\sim p(\mathbf{x}_k | \mathbf{x}_{k-1}) \\ \mathbf{y}_k &\sim p(\mathbf{y}_k | \mathbf{x}_k), \end{aligned} \quad (3.111)$$

where $\mathbf{x}_k \in \mathbb{R}^n$ is the state on time step k and $\mathbf{y}_k \in \mathbb{R}^m$ is the measurement. The state and measurements may contain both discrete and continuous components.

The SIS algorithm uses a weighted set of particles $\{(w_k^{(i)}, \mathbf{x}_k^{(i)}) : i = 1, \dots, N\}$ for representing the filtering distribution $p(\mathbf{x}_k | \mathbf{y}_{1:k})$ such that on every time step k an approximation of the expectation of an arbitrary function $\mathbf{g}(\mathbf{x})$ can be calculated as the weighted sample average

$$E[\mathbf{g}(\mathbf{x}_k) | \mathbf{y}_{1:k}] \approx \sum_{i=1}^N w_k^{(i)} \mathbf{g}(\mathbf{x}_k^{(i)}). \quad (3.112)$$

Equivalently, SIS can be interpreted to form an approximation of the posterior distribution as

$$p(\mathbf{x}_k | \mathbf{y}_{1:k}) \approx \sum_{i=1}^N w_k^{(i)} \delta(\mathbf{x}_k - \mathbf{x}_k^{(i)}), \quad (3.113)$$

where $\delta(\cdot)$ is the Dirac delta function.

The generic sequential importance sampling algorithm can be now described as follows:

Algorithm 3.15 (Sequential importance sampling). *Steps of SIS are the following:*

1. **Initialization:** Draw N samples $\mathbf{x}_0^{(i)}$ from the prior

$$\mathbf{x}_0^{(i)} \sim p(\mathbf{x}_0) \quad (3.114)$$

and set

$$w_0^{(i)} = 1/N \quad (3.115)$$

2. **Prediction:** Draw N new samples $\mathbf{x}_k^{(i)}$ from importance distribution

$$\mathbf{x}_k^{(i)} \sim \pi(\mathbf{x}_k | \mathbf{x}_{0:k-1}^{(i)}, \mathbf{y}_{1:k}) \quad (3.116)$$

3. **Update:** Calculate new weights according to

$$w_k^{(i)} = w_{k-1}^{(i)} \frac{p(\mathbf{y}_k | \mathbf{x}_k^{(i)}) p(\mathbf{x}_k^{(i)} | \mathbf{x}_{k-1}^{(i)})}{\pi(\mathbf{x}_k^{(i)} | \mathbf{x}_{0:k-1}^{(i)}, \mathbf{y}_{1:k})} \quad (3.117)$$

and normalize them to sum to unity.

4. Set $k \leftarrow k + 1$ and go to step 2.

3.4.2 Sequential Importance Resampling

One problem in the SIS algorithm described in the previous section is that we very easily encounter the situation that almost all the particles have zero weights and only a few of them (or only one) are non-zero. This is called the *degeneracy* problem in particle filtering literature and it used to prevent practical applications of particle filters for long time.

The degeneracy problem can be solved by using *resampling* procedure. It refers to a procedure where we draw N new samples from the discrete distribution defined by the weights and replace the old set of N samples with this new set. This procedure can be written as the following algorithm:

Algorithm 3.16 (Resampling). *Resampling procedure can be described as follows:*

1. Interpret each weight $w_k^{(i)}$ as the probability of obtaining the sample index i in the set $\{\mathbf{x}_k^{(i)} \mid i = 1, \dots, N\}$.
2. Draw N samples from that discrete distribution and replace the old sample set with this new one.
3. Set all weights to the constant value $w_k^{(i)} = 1/N$.

The idea of the resampling procedure is to remove particles with very small weights and duplicate particles with large weights. Although, the theoretical distribution represented by the weighted set of samples does not change, resampling induces additional variance to estimates. This variance introduced by the resampling procedure can be reduced by proper choice of the resampling method. The *stratified resampling* algorithm (Kitagawa, 1996) is optimal in terms of variance.

*Sequential importance resampling (SIR)*² (Gordon et al., 1993; Kitagawa, 1996; Doucet et al., 2001; Ristic et al., 2004), is a generalization of the *particle filtering* framework, in which the resampling step is included as part of the sequential importance sampling algorithm.

Usually the resampling is not performed on every time step, but only when it is actually needed. One way of implementing this is to do resampling on every n th step, where n is some predefined constant. This method has the advantage that it is unbiased. Another way, which is used here, is the *adaptive resampling*. In this method, the effective number of particles, which is estimated from the variance of the particle weights (Liu and Chen, 1995), is used for monitoring the need for resampling. The estimate for the effective number of particles can be computed as:

$$n_{\text{eff}} \approx \frac{1}{\sum_{i=1}^N \left(w_k^{(i)}\right)^2}, \quad (3.118)$$

where $w_k^{(i)}$ is the normalized weight of particle i on the time step k (Liu and Chen, 1995). Resampling is performed when the effective number of particles is significantly less than the total number of particles, for example, $n_{\text{eff}} < N/10$, where N is the total number of particles.

Algorithm 3.17 (Sequential importance resampling). *The SIR algorithm can be summarized as follows:*

1. Draw new point $\mathbf{x}_k^{(i)}$ for each point in the sample set $\{\mathbf{x}_{k-1}^{(i)}, i = 1, \dots, N\}$ from the importance distribution:

$$\mathbf{x}_k^{(i)} \sim \pi(\mathbf{x}_k \mid \mathbf{x}_{k-1}^{(i)}, \mathbf{y}_{1:k}), \quad i = 1, \dots, N. \quad (3.119)$$

²Sequential importance resampling (SIR) is also often referred to as *sampling importance resampling (SIR)* or *sequential importance sampling resampling (SISR)*.

2. Calculate new weights

$$w_k^{(i)} = w_{k-1}^{(i)} \frac{p(\mathbf{y}_k | \mathbf{x}_k^{(i)}) p(\mathbf{x}_k^{(i)} | \mathbf{x}_{k-1}^{(i)})}{\pi(\mathbf{x}_k^{(i)} | \mathbf{x}_{k-1}^{(i)}, \mathbf{y}_{1:k})}, \quad i = 1, \dots, N. \quad (3.120)$$

and normalize them to sum to unity.

3. If the effective number of particles (3.118) is too low, perform resampling.

The performance of the SIR algorithm is depends on the quality of the importance distribution $\pi(\cdot)$, which is an approximation to posterior distribution of states given the values at the previous step. The importance distribution should be in such functional form that we can easily draw samples from it and that it is possible to evaluate the probability densities of the sample points. *The optimal importance distribution* in terms of variance (see, e.g., Doucet et al., 2001; Ristic et al., 2004) is

$$\pi(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{1:k}) = p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{y}_{1:k}). \quad (3.121)$$

If the optimal importance distribution cannot be directly used, good importance distributions can be obtained by *local linearization* where a mixture of extended Kalman filters (EKF) or unscented Kalman filters (UKF) is used as the importance distribution (Doucet et al., 2000; Van der Merwe et al., 2001). Van der Merwe et al. (2001) also suggest a Metropolis-Hastings step after (or in place of) resampling step to smooth the resulting distribution, but from their results, it seems that this extra computation step has no significant performance effect. A particle filter with UKF importance distribution is also referred to as *unscented particle filter* (UPF).

By tuning the resampling algorithm to specific estimation problems and possibly changing the order of weight computation and sampling, accuracy and computational efficiency of the algorithm can be improved (Fearnhead and Clifford, 2003). An important issue is that sampling is more efficient without replacement, such that duplicate samples are not stored. There is also evidence that in some situations it is more efficient to use a simple deterministic algorithm for preserving the N most likely particles. In the article (Punskaya et al., 2002) it is shown that in digital demodulation, where the sampled space is discrete and the optimization criterion is the minimum error, the deterministic algorithm performs better.

The bootstrap filter (Gordon et al., 1993) is a variation of SIR, where the dynamic model $p(\mathbf{x}_k | \mathbf{x}_{k-1})$ is used as the importance distribution. This makes the implementation of the algorithm very easy, but due to the inefficiency of the importance distribution it may require a very large number of Monte Carlo samples for accurate estimation results. In bootstrap filter the resampling is normally done at each time step.

Algorithm 3.18 (Bootstrap filter). *The bootstrap filter algorithm is given as follows:*

1. Draw new point $\mathbf{x}_k^{(i)}$ for each point in the sample set $\{\mathbf{x}_{k-1}^{(i)}, i = 1, \dots, N\}$ from the dynamic model:

$$\mathbf{x}_k^{(i)} \sim p(\mathbf{x}_k | \mathbf{x}_{k-1}^{(i)}), \quad i = 1, \dots, N. \quad (3.122)$$

2. Calculate the weights

$$w_k^{(i)} = p(\mathbf{y}_k | \mathbf{x}_k^{(i)}), \quad i = 1, \dots, N. \quad (3.123)$$

and normalize them to sum to unity.

3. Do resampling.

Another variation of sequential importance resampling is the auxiliary SIR (ASIR) filter (Pitt and Shephard, 1999). The idea of the ASIR is to mimic the availability of optimal importance distribution by performing the resampling at step $k - 1$ using the available measurement at time k .

One problem encountered in particle filtering, despite the usage of resampling procedure, is called *sample impoverishment* (see, e.g., Ristic et al., 2004). It refers to the effect that when the noise in the dynamic model is very small, many of the particles in the particle set will turn out to have exactly the same value. That is, the resampling step simply multiplies a few (or one) particles and thus we end up having a set of identical copies of certain high weighted particles. This problem can be diminished by using, for example, resample-move algorithm, regularization or MCMC steps (Ristic et al., 2004).

Because low noise in the dynamic model causes problems with the sample impoverishment, it also implies that pure recursive estimation with particle filters is challenging. This is because in pure recursive estimation the process noise is formally zero and thus a basic SIR based particle filter is likely to perform very badly. However, pure recursive estimation, such as recursive estimation of static parameters can be done by applying a Rao-Blackwellized particle filter instead of a basic SIR particle filter.

3.4.3 Rao-Blackwellized Particle Filter

One way of improving the efficiency of SIR is to use Rao-Blackwellization. The idea of the *Rao-Blackwellized particle filter* (RBPF) (Akashi and Kumamoto, 1977; Doucet et al., 2001; Ristic et al., 2004) is that sometimes it is possible to evaluate some of the filtering equations analytically and the others with Monte Carlo sampling instead of computing everything with pure sampling. According to the *Rao-Blackwell theorem* (see, e.g., Berger, 1985; Casella and Robert, 1996) this leads to estimators with less variance than what could be obtained with pure Monte Carlo sampling. An intuitive way of understanding this is that the marginalization replaces the finite Monte Carlo particle set representation with an infinite closed form particle set, which is always more accurate than any finite set.

Most commonly Rao-Blackwellized particle filtering refers to marginalized filtering of conditionally Gaussian Markov models of the form

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{x}_{k-1}, \boldsymbol{\theta}_{k-1}) &= \mathcal{N}(\mathbf{x}_k | \mathbf{A}_{k-1}(\boldsymbol{\theta}_{k-1}) \mathbf{x}_{k-1}, \mathbf{Q}_{k-1}(\boldsymbol{\theta}_{k-1})) \\ p(\mathbf{y}_k | \mathbf{x}_k, \boldsymbol{\theta}_k) &= \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k(\boldsymbol{\theta}_k) \mathbf{x}_k, \mathbf{R}_k(\boldsymbol{\theta}_k)) \\ p(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{k-1}) &= (\text{any given form}), \end{aligned} \quad (3.124)$$

where \mathbf{x}_k is the state, \mathbf{y}_k is the measurement, and $\boldsymbol{\theta}_k$ is an arbitrary latent variable. If also the prior of \mathbf{x}_k is Gaussian, due to conditionally Gaussian structure of the model the state variables \mathbf{x}_k can be integrated out analytically and only the latent variables $\boldsymbol{\theta}_k$ need to be sampled. The Rao-Blackwellized particle filter uses SIR for the latent variables and computes everything else in closed form.

Algorithm 3.19 (Conditionally Gaussian Rao-Blackwellized particle filter). *Given an importance distribution $\pi(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{1:k-1}^{(i)}, \mathbf{y}_{1:k})$ and a set of weighted samples $\{w_{k-1}^{(i)}, \boldsymbol{\theta}_{k-1}^{(i)}, \mathbf{m}_{k-1}^{(i)}, \mathbf{P}_{k-1}^{(i)} : i = 1, \dots, N\}$, the Rao-Blackwellized particle filter processes each measurement \mathbf{y}_k as follows (Doucet et al., 2001):*

1. Perform Kalman filter predictions for each of the Kalman filter means and covariances in the particles $i = 1, \dots, N$ conditional on the previously drawn latent variable values $\boldsymbol{\theta}_{k-1}^{(i)}$

$$\begin{aligned} \mathbf{m}_k^{-(i)} &= \mathbf{A}_{k-1}(\boldsymbol{\theta}_{k-1}^{(i)}) \mathbf{m}_{k-1}^{(i)} \\ \mathbf{P}_k^{-(i)} &= \mathbf{A}_{k-1}(\boldsymbol{\theta}_{k-1}^{(i)}) \mathbf{P}_{k-1}^{(i)} \mathbf{A}_{k-1}^T(\boldsymbol{\theta}_{k-1}^{(i)}) + \mathbf{Q}_{k-1}(\boldsymbol{\theta}_{k-1}^{(i)}). \end{aligned} \quad (3.125)$$

2. Draw new latent variables $\boldsymbol{\theta}_k^{(i)}$ for each particle in $i = 1, \dots, N$ from the corresponding importance distributions

$$\boldsymbol{\theta}_k^{(i)} \sim \pi(\boldsymbol{\theta}_k | \boldsymbol{\theta}_{1:k-1}^{(i)}, \mathbf{y}_{1:k}). \quad (3.126)$$

3. Calculate new weights as follows:

$$w_k^{(i)} \propto w_{k-1}^{(i)} \frac{p(\mathbf{y}_k | \boldsymbol{\theta}_{1:k}^{(i)}, \mathbf{y}_{1:k-1}) p(\boldsymbol{\theta}_k^{(i)} | \boldsymbol{\theta}_{k-1}^{(i)})}{\pi(\boldsymbol{\theta}_k^{(i)} | \boldsymbol{\theta}_{1:k-1}^{(i)}, \mathbf{y}_{1:k})}, \quad (3.127)$$

where the likelihood term is the marginal measurement likelihood of the Kalman filter

$$\begin{aligned} p(\mathbf{y}_k | \boldsymbol{\theta}_{1:k}^{(i)}, \mathbf{y}_{1:k-1}) \\ = \mathcal{N}\left(\mathbf{y}_k \mid \mathbf{H}_k(\boldsymbol{\theta}_k^{(i)}) \mathbf{m}_k^{-(i)}, \mathbf{H}_k(\boldsymbol{\theta}_k^{(i)}) \mathbf{P}_k^{-(i)} \mathbf{H}_k^T(\boldsymbol{\theta}_k^{(i)}) + \mathbf{R}_k(\boldsymbol{\theta}_k^{(i)})\right). \end{aligned} \quad (3.128)$$

such that the model parameters in the Kalman filter are conditioned on the drawn latent variable value $\boldsymbol{\theta}_k^{(i)}$. Then normalize the weights to sum to unity.

4. Perform Kalman filter updates for each of the particles conditional on the drawn latent variables $\boldsymbol{\theta}_k^{(i)}$

$$\begin{aligned}
\mathbf{v}_k^{(i)} &= \mathbf{y}_k - \mathbf{H}_k(\boldsymbol{\theta}_k^{(i)}) \mathbf{m}_k^- \\
\mathbf{S}_k^{(i)} &= \mathbf{H}_k(\boldsymbol{\theta}_k^{(i)}) \mathbf{P}_k^{- (i)} \mathbf{H}_k^T(\boldsymbol{\theta}_k^{(i)}) + \mathbf{R}_k(\boldsymbol{\theta}_k^{(i)}) \\
\mathbf{K}_k^{(i)} &= \mathbf{P}_k^{- (i)} \mathbf{H}_k^T(\boldsymbol{\theta}_k^{(i)}) \mathbf{S}_k^{-1} \\
\mathbf{m}_k^{(i)} &= \mathbf{m}_k^{- (i)} + \mathbf{K}_k^{(i)} \mathbf{v}_k^{(i)} \\
\mathbf{P}_k^{(i)} &= \mathbf{P}_k^{- (i)} - \mathbf{K}_k^{(i)} \mathbf{S}_k^{(i)} [\mathbf{K}_k^{(i)}]^T.
\end{aligned} \tag{3.129}$$

5. If the effective number of particles (3.118) is too low, perform resampling.

The Rao-Blackwellized particle filter produces for each time step k a set of weighted samples $\{w_k^{(i)}, \boldsymbol{\theta}_k^{(i)}, \mathbf{m}_k^{(i)}, \mathbf{P}_k^{(i)} : i = 1, \dots, N\}$ such that expectation of a function $\mathbf{g}(\cdot)$ can be approximated as

$$\mathbb{E}[\mathbf{g}(\mathbf{x}_k, \boldsymbol{\theta}_k) | \mathbf{y}_{1:k}] \approx \sum_{i=1}^N w_k^{(i)} \int \mathbf{g}(\mathbf{x}_k, \boldsymbol{\theta}_k^{(i)}) \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^{(i)}, \mathbf{P}_k^{(i)}) d\mathbf{x}_k. \tag{3.130}$$

Equivalently the RBPF can be interpreted to form an approximation of the filtering distribution as

$$p(\mathbf{x}_k, \boldsymbol{\theta}_k | \mathbf{y}_{1:k}) \approx \sum_{i=1}^N w_k^{(i)} \delta(\boldsymbol{\theta}_k - \boldsymbol{\theta}_k^{(i)}) \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^{(i)}, \mathbf{P}_k^{(i)}). \tag{3.131}$$

In some cases, when the filtering model is not strictly Gaussian due to slight non-linearities in either dynamic or measurement models it is possible to replace the exact Kalman filter update and prediction steps in RBPF with extended Kalman filter (EKF) or unscented Kalman filter (UKF) prediction and update steps.

In addition to the conditional Gaussian models, another general class of models where Rao-Blackwellization can often be applied are state space models with unknown static parameters. These models are of the form (Storvik, 2002)

$$\begin{aligned}
\mathbf{x}_k &\sim p(\mathbf{x}_k | \mathbf{x}_{k-1}, \boldsymbol{\theta}) \\
\mathbf{y}_k &\sim p(\mathbf{y}_k | \mathbf{x}_k, \boldsymbol{\theta}) \\
\boldsymbol{\theta} &\sim p(\boldsymbol{\theta}),
\end{aligned} \tag{3.132}$$

where vector $\boldsymbol{\theta}$ contains the unknown static parameters. If the posterior distribution of parameters $\boldsymbol{\theta}$ depends only on some sufficient statistics

$$\mathbf{T}_k = \mathbf{T}_k(\mathbf{x}_{1:k}, \mathbf{y}_{1:k}), \tag{3.133}$$

and if the sufficient statics are easy to update recursively, then sampling of the state and parameters can be efficiently performed by recursively computing the sufficient statistics conditionally to the sampled states and the measurements (Storvik, 2002).

A particularly useful special case is obtained when the dynamic model is independent of the parameters θ . In this case, if conditionally to the state \mathbf{x}_k the prior $p(\theta)$ belongs to the conjugate family of the likelihood $p(\mathbf{y}_k | \mathbf{x}_k, \theta)$, the static parameters θ can be marginalized out and only the states need to be sampled.

Chapter 4

Optimal Smoothing

4.1 Formal Equations and Exact Solutions

4.1.1 Optimal Smoothing Equations

The purpose of *optimal smoothing*¹ is to compute the marginal posterior distribution of the state \mathbf{x}_k at the time step k after receiving the measurements up to a time step T , where $T > k$:

$$p(\mathbf{x}_k | \mathbf{y}_{1:T}). \quad (4.1)$$

The difference between filters and smoothers is that *the optimal filter* computes its estimates using only the measurements obtained before and on the time step k , but *the optimal smoother* uses also the future measurements for computing its estimates. After obtaining the filtering posterior state distributions, the following theorem gives the equations for computing the marginal posterior distributions for each time step conditionally to all measurements up to the time step T :

Theorem 4.1 (Bayesian optimal smoothing equations). *The backward recursive equations for computing the smoothed distributions $p(\mathbf{x}_k | \mathbf{y}_{1:T})$ for any $k < T$ are given by the following Bayesian (fixed interval) smoothing equations*

$$\begin{aligned} p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k}) &= \int p(\mathbf{x}_{k+1} | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k}) d\mathbf{x}_k \\ p(\mathbf{x}_k | \mathbf{y}_{1:T}) &= p(\mathbf{x}_k | \mathbf{y}_{1:k}) \int \left[\frac{p(\mathbf{x}_{k+1} | \mathbf{x}_k) p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T})}{p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})} \right] d\mathbf{x}_{k+1}, \end{aligned} \quad (4.2)$$

where $p(\mathbf{x}_k | \mathbf{y}_{1:k})$ is the filtering distribution of the time step k . Note that the term $p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})$ is simply the predicted distribution of time step $k + 1$. The integrations are replaced by summations if some of the state components are discrete.

Proof. Due to the Markov properties the state \mathbf{x}_k is independent of $\mathbf{y}_{k+1:T}$ given \mathbf{x}_{k+1} , which gives $p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) = p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:k})$. By using the Bayes'

¹In this document only fixed-interval smoothing is considered.

rule the distribution of \mathbf{x}_k given \mathbf{x}_{k+1} and $\mathbf{y}_{1:T}$ can be expressed as

$$\begin{aligned}
 p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) &= p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:k}) \\
 &= \frac{p(\mathbf{x}_k, \mathbf{x}_{k+1} | \mathbf{y}_{1:k})}{p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})} \\
 &= \frac{p(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{y}_{1:k}) p(\mathbf{x}_k | \mathbf{y}_{1:k})}{p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})} \\
 &= \frac{p(\mathbf{x}_{k+1} | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k})}{p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})}.
 \end{aligned} \tag{4.3}$$

The joint distribution of \mathbf{x}_k and \mathbf{x}_{k+1} given $\mathbf{y}_{1:T}$ can be now computed as

$$\begin{aligned}
 p(\mathbf{x}_k, \mathbf{x}_{k+1} | \mathbf{y}_{1:T}) &= p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T}) \\
 &= p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:k}) p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T}) \\
 &= \frac{p(\mathbf{x}_{k+1} | \mathbf{x}_k) p(\mathbf{x}_k | \mathbf{y}_{1:k}) p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T})}{p(\mathbf{x}_{k+1} | \mathbf{y}_{1:k})},
 \end{aligned} \tag{4.4}$$

where $p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T})$ is the smoothed distribution of the time step $k + 1$. The marginal distribution of \mathbf{x}_k given $\mathbf{y}_{1:T}$ is given by integral (or summation) over \mathbf{x}_{k+1} in Equation (4.4), which gives the desired result. \square

4.1.2 Discrete-Time Rauch-Tung-Striebel Smoother

The *discrete-time Rauch-Tung-Striebel (RTS)*² (see, e.g., Rauch et al., 1965; Gelb, 1974; Bar-Shalom et al., 2001) can be used for computing the closed form smoothing solution

$$p(\mathbf{x}_k | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^s, \mathbf{P}_k^s), \tag{4.5}$$

to the linear filtering model (3.17). The difference to the solution computed by the *Kalman filter* is that the smoothed solution is conditional on the whole measurement data $\mathbf{y}_{1:T}$, while the filtering solution is conditional only on the measurements obtained before and on the time step k , that is, on the measurements $\mathbf{y}_{1:k}$.

Theorem 4.2 (Discrete-time RTS smoother). *The backward recursion equations for the discrete-time fixed interval Rauch-Tung-Striebel smoother (Kalman smoother) are given as*

$$\begin{aligned}
 \mathbf{m}_{k+1}^- &= \mathbf{A}_k \mathbf{m}_k \\
 \mathbf{P}_{k+1}^- &= \mathbf{A}_k \mathbf{P}_k \mathbf{A}_k^T + \mathbf{Q}_k \\
 \mathbf{C}_k &= \mathbf{P}_k \mathbf{A}_k^T [\mathbf{P}_{k+1}^-]^{-1} \\
 \mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k [\mathbf{m}_{k+1}^s - \mathbf{m}_{k+1}^-] \\
 \mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k [\mathbf{P}_{k+1}^s - \mathbf{P}_{k+1}^-] \mathbf{C}_k^T,
 \end{aligned} \tag{4.6}$$

²Also called discrete-time Kalman smoother.

where \mathbf{m}_k and \mathbf{P}_k are the mean and covariance computed by the Kalman filter. The recursion is started from the last time step T , with $\mathbf{m}_T^s = \mathbf{m}_T$ and $\mathbf{P}_T^s = \mathbf{P}_T$. Note that the first two of the equations are simply the Kalman filter prediction equations.

Proof. Similarly to the Kalman filter case, by Lemma A.1, the joint distribution of \mathbf{x}_k and \mathbf{x}_{k+1} given $\mathbf{y}_{1:k}$ is

$$\begin{aligned} p(\mathbf{x}_k, \mathbf{x}_{k+1} \mid \mathbf{y}_{1:k}) &= p(\mathbf{x}_{k+1} \mid \mathbf{x}_k) p(\mathbf{x}_k \mid \mathbf{y}_{1:k}) \\ &= \mathcal{N}(\mathbf{x}_{k+1} \mid \mathbf{A}_k \mathbf{x}_k, \mathbf{Q}_k) \mathcal{N}(\mathbf{x}_k \mid \mathbf{m}_k, \mathbf{P}_k) \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_{k+1} \end{bmatrix} \mid \mathbf{m}_1, \mathbf{P}_1\right), \end{aligned} \quad (4.7)$$

where

$$\mathbf{m}_1 = \begin{pmatrix} \mathbf{m}_k \\ \mathbf{A}_k \mathbf{m}_k \end{pmatrix}, \quad \mathbf{P}_1 = \begin{pmatrix} \mathbf{P}_k & \mathbf{P}_k \mathbf{A}_k^T \\ \mathbf{A}_k \mathbf{P}_k & \mathbf{A}_k \mathbf{P}_k \mathbf{A}_k^T + \mathbf{Q}_k \end{pmatrix}. \quad (4.8)$$

Due to the Markov property of the states we have

$$p(\mathbf{x}_k \mid \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) = p(\mathbf{x}_k \mid \mathbf{x}_{k+1}, \mathbf{y}_{1:k}), \quad (4.9)$$

and thus by Lemma A.2 we get the conditional distribution

$$\begin{aligned} p(\mathbf{x}_k \mid \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) &= p(\mathbf{x}_k \mid \mathbf{x}_{k+1}, \mathbf{y}_{1:k}) \\ &= \mathcal{N}(\mathbf{x}_k \mid \mathbf{m}_2, \mathbf{P}_2), \end{aligned} \quad (4.10)$$

where

$$\begin{aligned} \mathbf{C}_k &= \mathbf{P}_k \mathbf{A}_k^T (\mathbf{A}_k \mathbf{P}_k \mathbf{A}_k^T + \mathbf{Q}_k)^{-1} \\ \mathbf{m}_2 &= \mathbf{m}_k + \mathbf{C}_k (\mathbf{x}_{k+1} - \mathbf{A}_k \mathbf{m}_k) \\ \mathbf{P}_2 &= \mathbf{P}_k - \mathbf{C}_k (\mathbf{A}_k \mathbf{P}_k \mathbf{A}_k^T + \mathbf{Q}_k) \mathbf{C}_k^T. \end{aligned} \quad (4.11)$$

The joint distribution of \mathbf{x}_k and \mathbf{x}_{k+1} given all the data is

$$\begin{aligned} p(\mathbf{x}_{k+1}, \mathbf{x}_k \mid \mathbf{y}_{1:T}) &= p(\mathbf{x}_k \mid \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) p(\mathbf{x}_{k+1} \mid \mathbf{y}_{1:T}) \\ &= \mathcal{N}(\mathbf{x}_k \mid \mathbf{m}_2, \mathbf{P}_2) \mathcal{N}(\mathbf{x}_{k+1} \mid \mathbf{m}_{k+1}^s, \mathbf{P}_{k+1}^s) \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_{k+1} \\ \mathbf{x}_k \end{bmatrix} \mid \mathbf{m}_3, \mathbf{P}_3\right) \end{aligned} \quad (4.12)$$

where

$$\begin{aligned} \mathbf{m}_3 &= \begin{pmatrix} \mathbf{m}_{k+1}^s \\ \mathbf{m}_k + \mathbf{C}_k (\mathbf{m}_{k+1}^s - \mathbf{A}_k \mathbf{m}_k) \end{pmatrix} \\ \mathbf{P}_3 &= \begin{pmatrix} \mathbf{P}_{k+1}^s & \mathbf{P}_{k+1}^s \mathbf{C}_k^T \\ \mathbf{C}_k \mathbf{P}_{k+1}^s & \mathbf{C}_k \mathbf{P}_{k+1}^s \mathbf{C}_k^T + \mathbf{P}_2 \end{pmatrix}. \end{aligned} \quad (4.13)$$

Thus by Lemma A.2, the marginal distribution of \mathbf{x}_k is given as

$$p(\mathbf{x}_k | \mathbf{y}_{1:T}) = N(\mathbf{x}_k | \mathbf{m}_k^s, \mathbf{P}_k^s), \quad (4.14)$$

where

$$\begin{aligned} \mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k (\mathbf{m}_{k+1}^s - \mathbf{A}_k \mathbf{m}_k) \\ \mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k (\mathbf{P}_{k+1}^s - \mathbf{A}_k \mathbf{P}_k \mathbf{A}_k^T - \mathbf{Q}_k) \mathbf{C}_k^T. \end{aligned} \quad (4.15)$$

□

Example 4.1 (RTS smoother for Gaussian random walk). *The RTS smoother for the random walk model given in Example 3.1 is given by the equations*

$$\begin{aligned} m_{k+1}^- &= m_k \\ P_{k+1}^- &= P_k + q \\ m_k^s &= m_k + \frac{P_k}{P_{k+1}^-} (m_{k+1}^s - m_{k+1}^-) \\ P_k^s &= P_k + \left(\frac{P_k}{P_{k+1}^-} \right)^2 [P_{k+1}^s - P_{k+1}^-], \end{aligned} \quad (4.16)$$

where m_k and P_k are the updated mean and covariance from the Kalman filter in Example 3.2.

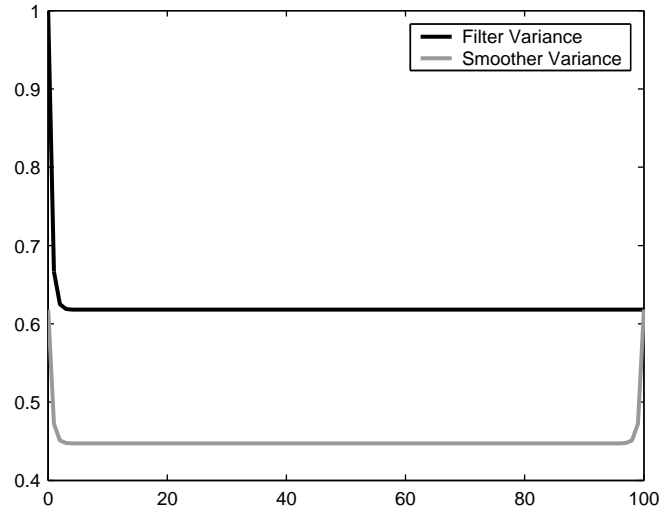


Figure 4.1: Filter and smoother variances in the Kalman smoothing example (Example 4.1).

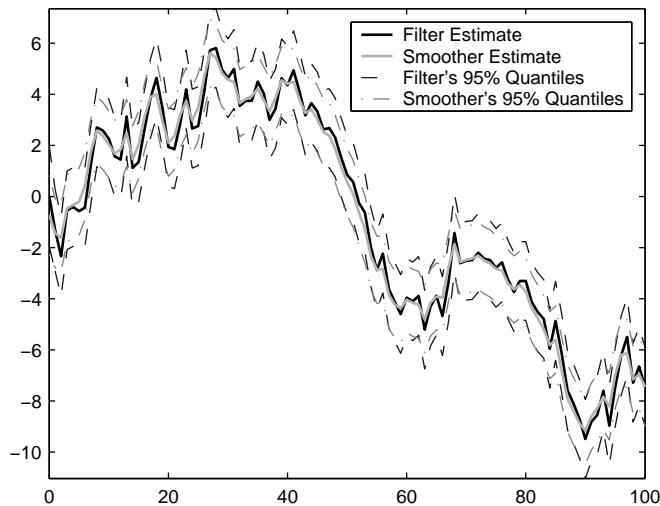


Figure 4.2: Filter and smoother estimates in the Kalman smoothing example (Example 4.1).

4.2 Gaussian Approximation Based Smoothing

4.2.1 Discrete-Time Extended Rauch-Tung-Striebel Smoother

The first order (i.e., linearized) extended Rauch-Tung-Striebel smoother (ERTSS) (Cox, 1964; Sage and Melsa, 1971) can be obtained from the basic RTS smoother equations by replacing the prediction equations with first order approximations. Higher order extended Kalman smoothers are also possible (see, e.g., Cox, 1964; Sage and Melsa, 1971), but only the first order version is presented here.

For the additive model Equation (3.52) the extended Rauch-Tung-Striebel smoother algorithm is the following:

Algorithm 4.1 (Extended RTS smoother). *The equations for the extended RTS smoother are*

$$\begin{aligned}
 \mathbf{m}_{k+1}^- &= \mathbf{f}(\mathbf{m}_k) \\
 \mathbf{P}_{k+1}^- &= \mathbf{F}_x(\mathbf{m}_k) \mathbf{P}_k \mathbf{F}_x^T(\mathbf{m}_k) + \mathbf{Q}_k \\
 \mathbf{C}_k &= \mathbf{P}_k \mathbf{F}_x^T(\mathbf{m}_k) [\mathbf{P}_{k+1}^-]^{-1} \\
 \mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k [\mathbf{m}_{k+1}^s - \mathbf{m}_{k+1}^-] \\
 \mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k [\mathbf{P}_{k+1}^s - \mathbf{P}_{k+1}^-] \mathbf{C}_k^T,
 \end{aligned} \tag{4.17}$$

where the matrix $\mathbf{F}_x(\mathbf{m}_k)$ is the Jacobian matrix of $\mathbf{f}(\mathbf{x})$ evaluated at \mathbf{m}_k .

The above procedure is a recursion, which can be used for computing the smoothing distribution of step k from the smoothing distribution of time step $k+1$.

Because the smoothing distribution and filtering distribution of the last time step T are the same, we have $\mathbf{m}_T^s = \mathbf{m}_T$, $\mathbf{P}_T^s = \mathbf{P}_T$, and thus the recursion can be used for computing the smoothing distributions of all time steps by starting from the last step $k = T$ and proceeding backwards to the initial step $k = 0$.

Proof. Assume that the approximate means and covariances of the filtering distributions

$$p(\mathbf{x}_k | \mathbf{y}_{1:k}) \approx N(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k),$$

for the model (3.52) have been computed by the extended Kalman filter or a similar method. Further assume that the smoothing distribution of time step $k + 1$ is known and approximately Gaussian

$$p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T}) \approx N(\mathbf{x}_{k+1} | \mathbf{m}_{k+1}^s, \mathbf{P}_{k+1}^s).$$

As in the derivation of the prediction step of EKF in Section 3.2.2, the approximate joint distribution of \mathbf{x}_k and \mathbf{x}_{k+1} given $\mathbf{y}_{1:k}$ is

$$p(\mathbf{x}_k, \mathbf{x}_{k+1} | \mathbf{y}_{1:k}) = N\left(\begin{bmatrix} \mathbf{x}_k \\ \mathbf{x}_{k+1} \end{bmatrix} \middle| \mathbf{m}_1, \mathbf{P}_1\right), \quad (4.18)$$

where

$$\begin{aligned} \mathbf{m}_1 &= \begin{pmatrix} \mathbf{m}_k \\ \mathbf{f}(\mathbf{m}_k) \end{pmatrix} \\ \mathbf{P}_1 &= \begin{pmatrix} \mathbf{P}_k & \mathbf{P}_k \mathbf{F}_x^T \\ \mathbf{F}_x \mathbf{P}_k & \mathbf{F}_x \mathbf{P}_k \mathbf{F}_x^T + \mathbf{Q}_k \end{pmatrix}. \end{aligned} \quad (4.19)$$

where the Jacobian matrix \mathbf{F}_x of $\mathbf{f}(\mathbf{x})$ is evaluated at $\mathbf{x} = \mathbf{m}_k$. By conditioning to \mathbf{x}_{k+1} as in RTS derivation in Section 4.1.2 we get

$$\begin{aligned} p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) &= p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:k}) \\ &= N(\mathbf{x}_k | \mathbf{m}_2, \mathbf{P}_2), \end{aligned} \quad (4.20)$$

where

$$\begin{aligned} \mathbf{C}_k &= \mathbf{P}_k \mathbf{F}_x^T (\mathbf{F}_x \mathbf{P}_k \mathbf{F}_x^T + \mathbf{Q}_k)^{-1} \\ \mathbf{m}_2 &= \mathbf{m}_k + \mathbf{C}_k (\mathbf{x}_{k+1} - \mathbf{f}(\mathbf{m}_k)) \\ \mathbf{P}_2 &= \mathbf{P}_k - \mathbf{C}_k (\mathbf{F}_x \mathbf{P}_k \mathbf{F}_x^T + \mathbf{Q}_k) \mathbf{C}_k^T. \end{aligned} \quad (4.21)$$

The joint distribution of \mathbf{x}_k and \mathbf{x}_{k+1} given all the data is now

$$\begin{aligned} p(\mathbf{x}_{k+1}, \mathbf{x}_k | \mathbf{y}_{1:T}) &= p(\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:T}) p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T}) \\ &= N\left(\begin{bmatrix} \mathbf{x}_{k+1} \\ \mathbf{x}_k \end{bmatrix} \middle| \mathbf{m}_3, \mathbf{P}_3\right) \end{aligned} \quad (4.22)$$

where

$$\begin{aligned} \mathbf{m}_3 &= \begin{pmatrix} \mathbf{m}_{k+1}^s \\ \mathbf{m}_k + \mathbf{C}_k (\mathbf{m}_{k+1}^s - \mathbf{f}(\mathbf{m}_k)) \end{pmatrix} \\ \mathbf{P}_3 &= \begin{pmatrix} \mathbf{P}_{k+1}^s & \mathbf{P}_{k+1}^s \mathbf{C}_k^T \\ \mathbf{C}_k \mathbf{P}_{k+1}^s & \mathbf{C}_k \mathbf{P}_{k+1}^s \mathbf{C}_k^T + \mathbf{P}_2 \end{pmatrix}. \end{aligned} \quad (4.23)$$

The marginal distribution of \mathbf{x}_k is then

$$p(\mathbf{x}_k | \mathbf{y}_{1:T}) = \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k^s, \mathbf{P}_k^s), \quad (4.24)$$

where

$$\begin{aligned} \mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k (\mathbf{m}_{k+1}^s - \mathbf{f}(\mathbf{m}_k)) \\ \mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k (\mathbf{P}_{k+1}^s - \mathbf{F}_x \mathbf{P}_k \mathbf{F}_x^T - \mathbf{Q}_k) \mathbf{C}_k^T. \end{aligned} \quad (4.25)$$

□

The generalization to non-additive model (3.65) is analogous to the filtering case.

4.2.2 Statistically Linearized RTS Smoother

The statistically linearized Rauch-Tung-Striebel smoother for the additive model (3.52) is the following:

Algorithm 4.2 (Statistically linearized RTS smoother). *The equations for the statistically linearized RTS smoother are*

$$\begin{aligned} \mathbf{m}_{k+1}^- &= \mathbb{E}[\mathbf{f}(\mathbf{x}_k)] \\ \mathbf{P}_{k+1}^- &= \mathbb{E}[\mathbf{f}(\mathbf{x}_k) \delta \mathbf{x}_k^T] \mathbf{P}_k^{-1} \mathbb{E}[\mathbf{f}(\mathbf{x}_k) \delta \mathbf{x}_k^T]^T + \mathbf{Q}_k \\ \mathbf{C}_k &= \mathbb{E}[\mathbf{f}(\mathbf{x}_k) \delta \mathbf{x}_k^T]^T [\mathbf{P}_{k+1}^-]^{-1} \\ \mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k [\mathbf{m}_{k+1}^s - \mathbf{m}_{k+1}^-] \\ \mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k [\mathbf{P}_{k+1}^s - \mathbf{P}_{k+1}^-] \mathbf{C}_k^T, \end{aligned} \quad (4.26)$$

where the expectations are taken with respect to the filtering distribution $\mathbf{x}_k \sim \mathcal{N}(\mathbf{m}_k, \mathbf{P}_k)$.

Proof. Analogous to the EKF case. □

The generalization to the non-additive case is also straight-forward.

4.2.3 Unscented Rauch-Tung-Striebel Smoother

The *unscented Rauch-Tung-Striebel smoother* (URTSS) (see, e.g., Särkkä, 2008) is a Gaussian approximation based smoother, where the non-linearity is approximated using the unscented transform. The smoother equations for the *non-additive model* (3.65) are given as follows:

Algorithm 4.3 (Unscented Rauch-Tung-Striebel smoother II). *A single step of the unscented RTS smoother is as follows:*

1. Form the matrix of sigma points of the n' -dimensional augmented random variable $\tilde{\mathbf{x}}_k = (\mathbf{x}_k^T \mathbf{q}_k^T)^T$

$$\tilde{\mathbf{X}}_k = (\tilde{\mathbf{m}}_k \quad \cdots \quad \tilde{\mathbf{m}}_k) + \sqrt{n' + \lambda} \begin{pmatrix} \mathbf{0} & \sqrt{\tilde{\mathbf{P}}_k} & -\sqrt{\tilde{\mathbf{P}}_k} \end{pmatrix}.$$

where

$$\tilde{\mathbf{m}}_k = \begin{pmatrix} \mathbf{m}_k \\ \mathbf{0} \end{pmatrix} \quad \tilde{\mathbf{P}}_k = \begin{pmatrix} \mathbf{P}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_k \end{pmatrix}.$$

2. Propagate the sigma points through the dynamic model:

$$\tilde{\mathbf{X}}_{k+1,i}^- = \mathbf{f}(\tilde{\mathbf{X}}_{k,i}^x, \tilde{\mathbf{X}}_{k,i}^q), \quad i = 1 \dots 2n' + 1,$$

where $\tilde{\mathbf{X}}_{k,i}^x$ and $\tilde{\mathbf{X}}_{k,i}^q$ denote the parts of the augmented sigma point i , which correspond to \mathbf{x}_k and \mathbf{q}_k , respectively.

3. Compute the predicted mean \mathbf{m}_{k+1}^- , the predicted covariance \mathbf{P}_{k+1}^- and the cross-covariance \mathbf{D}_{k+1} :

$$\begin{aligned} \mathbf{m}_{k+1}^- &= \sum_i W_{i-1}^{(m)} \tilde{\mathbf{X}}_{k+1,i}^- \\ \mathbf{P}_{k+1}^- &= \sum_i W_{i-1}^{(c)} (\tilde{\mathbf{X}}_{k+1,i}^- - \mathbf{m}_{k+1}^-) (\tilde{\mathbf{X}}_{k+1,i}^- - \mathbf{m}_{k+1}^-)^T \\ \mathbf{D}_{k+1} &= \sum_i W_{i-1}^{(c)} (\tilde{\mathbf{X}}_{k,i}^x - \mathbf{m}_k) (\tilde{\mathbf{X}}_{k+1,i}^- - \mathbf{m}_{k+1}^-)^T, \end{aligned} \quad (4.27)$$

where the definitions of the weights $W_i^{(m)}$ and $W_i^{(c)}$ are the same as in Section 3.2.5.

4. Compute the smoother gain \mathbf{C}_k , the smoothed mean \mathbf{m}_k^s and the covariance \mathbf{P}_k^s :

$$\begin{aligned} \mathbf{C}_k &= \mathbf{D}_{k+1} [\mathbf{P}_{k+1}^-]^{-1} \\ \mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k [\mathbf{m}_{k+1}^s - \mathbf{m}_{k+1}^-] \\ \mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k [\mathbf{P}_{k+1}^s - \mathbf{P}_{k+1}^-] \mathbf{C}_k^T. \end{aligned} \quad (4.28)$$

Proof. Assume that the approximate means and covariances of the filtering distributions are available:

$$p(\mathbf{x}_k | \mathbf{y}_{1:k}) \approx \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k),$$

and the smoothing distribution of time step $k + 1$ is known and approximately Gaussian

$$p(\mathbf{x}_{k+1} | \mathbf{y}_{1:T}) \approx \mathcal{N}(\mathbf{x}_{k+1} | \mathbf{m}_{k+1}^s, \mathbf{P}_{k+1}^s).$$

An unscented transform based approximation to the optimal smoothing solution can be derived as follows:

1. Generate unscented transform based Gaussian approximation to the joint distribution of \mathbf{x}_k and \mathbf{x}_{k+1} :

$$\begin{pmatrix} \mathbf{x}_k \\ \mathbf{x}_{k+1} \end{pmatrix} | \mathbf{y}_{1:k} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{m}_k \\ \mathbf{m}_{k+1}^- \end{pmatrix}, \begin{pmatrix} \mathbf{P}_k & \mathbf{D}_{k+1} \\ \mathbf{D}_{k+1}^T & \mathbf{P}_{k+1}^- \end{pmatrix} \right), \quad (4.29)$$

This can be done by concatenating the state and process noise to a new augmented random variable $\tilde{\mathbf{x}}_k = (\mathbf{x}_k^T \mathbf{q}_k^T)^T$, which then has the distribution

$$\tilde{\mathbf{x}}_k | \mathbf{y}_{1:k} \sim \mathcal{N} \left(\begin{pmatrix} \mathbf{m}_k \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{P}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_k \end{pmatrix} \right).$$

It is now easy to use the unscented transform for forming a Gaussian approximation to the joint distribution of $\tilde{\mathbf{x}}_k = (\mathbf{x}_k^T \mathbf{q}_k^T)^T$ and $\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{q}_k)$. The Gaussian approximation to the joint distribution of \mathbf{x}_k and \mathbf{x}_{k+1} can be formed by extracting the relevant parts of the mean and covariance from the joint Gaussian approximation of $\tilde{\mathbf{x}}_k$ and \mathbf{x}_{k+1} . This is done in Equations (4.27).

2. Because the distribution (4.29) is Gaussian, by the computation rules of Gaussian distributions and the conditional distribution of \mathbf{x}_k is given as

$$\mathbf{x}_k | \mathbf{x}_{k+1}, \mathbf{y}_{1:T} \sim \mathcal{N}(\mathbf{m}_2, \mathbf{P}_2),$$

where

$$\begin{aligned} \mathbf{C}_k &= \mathbf{D}_{k+1} [\mathbf{P}_{k+1}^-]^{-1} \\ \mathbf{m}_2 &= \mathbf{m}_k + \mathbf{C}_k (\mathbf{x}_{k+1} - \mathbf{m}_{k+1}^-) \\ \mathbf{P}_2 &= \mathbf{P}_k - \mathbf{C}_k \mathbf{P}_{k+1}^- \mathbf{C}_k^T. \end{aligned}$$

3. The rest of the derivation is completely analogous to the derivation of ERTSS in Section 4.2.1.

□

As the noises in the state space model (3.52) appear in additive manner, for that model it is possible write the URTSS equations in a bit simpler additive form:

Algorithm 4.4 (Unscented Rauch-Tung-Striebel smoother I). *The additive form unscented RTS smoother algorithm is the following:*

1. Form the matrix of sigma points:

$$\mathbf{X}_k = [\mathbf{m}_k \quad \cdots \quad \mathbf{m}_k] + \sqrt{n + \lambda} [\mathbf{0} \quad \sqrt{\mathbf{P}_k} \quad -\sqrt{\mathbf{P}_k}].$$

2. Propagate the sigma points through the dynamic model:

$$\hat{\mathbf{X}}_{k+1,i} = \mathbf{f}(\mathbf{X}_{k,i}), \quad i = 1 \dots 2n + 1.$$

3. Compute the predicted mean \mathbf{m}_{k+1}^- , the predicted covariance \mathbf{P}_{k+1}^- and the cross-covariance \mathbf{D}_{k+1} :

$$\begin{aligned} \mathbf{m}_{k+1}^- &= \sum_i W_{i-1}^{(m)} \hat{\mathbf{X}}_{k+1,i} \\ \mathbf{P}_{k+1}^- &= \sum_i W_{i-1}^{(c)} (\hat{\mathbf{X}}_{k+1,i} - \mathbf{m}_{k+1}^-) (\hat{\mathbf{X}}_{k+1,i} - \mathbf{m}_{k+1}^-)^T + \mathbf{Q}_k \\ \mathbf{D}_{k+1} &= \sum_i W_{i-1}^{(c)} (\mathbf{X}_{k,i} - \mathbf{m}_k) (\hat{\mathbf{X}}_{k+1,i} - \mathbf{m}_{k+1}^-)^T. \end{aligned} \quad (4.30)$$

4. Compute the smoother gain \mathbf{C}_k , the smoothed mean \mathbf{m}_k^s and the covariance \mathbf{P}_k^s as follows:

$$\begin{aligned} \mathbf{C}_k &= \mathbf{D}_{k+1} [\mathbf{P}_{k+1}^-]^{-1} \\ \mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k (\mathbf{m}_{k+1}^s - \mathbf{m}_{k+1}^-) \\ \mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k (\mathbf{P}_{k+1}^s - \mathbf{P}_{k+1}^-) \mathbf{C}_k^T. \end{aligned} \quad (4.31)$$

The above computations are started from the filtering result of the last time step $\mathbf{m}_T^s = \mathbf{m}_T$, $\mathbf{P}_T^s = \mathbf{P}_T$ and the recursion runs backwards for $k = T - 1, \dots, 0$.

4.2.4 Gaussian Assumed Density RTS Smoother

The Gaussian moment matching described in Section 3.2.7 can be used in smoothers in analogous manner as in Gaussian assumed density filters in Section 3.2.8. If we follow the extended RTS smoother derivation in Section 4.2.1, we get the following algorithm (see, e.g., Särkkä and Hartikainen, 2010a,b):

Algorithm 4.5 (Gaussian assumed density smoother). *The equations of the Gaussian assumed density RTS smoother are the following:*

$$\begin{aligned}
\mathbf{m}_{k+1}^- &= \int \mathbf{f}(\mathbf{x}_k) \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k) d\mathbf{x}_k \\
\mathbf{P}_{k+1|k}^- &= \int [\mathbf{f}(\mathbf{x}_k) - \mathbf{m}_{k+1}^-] [\mathbf{f}(\mathbf{x}_k) - \mathbf{m}_{k+1}^-]^T \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k) d\mathbf{x}_k + \mathbf{Q}_k \\
\mathbf{D}_{k+1} &= \int [\mathbf{x}_k - \mathbf{m}_k] [\mathbf{f}(\mathbf{x}_k) - \mathbf{m}_{k+1}^-]^T \mathcal{N}(\mathbf{x}_k | \mathbf{m}_k, \mathbf{P}_k) d\mathbf{x}_k \\
\mathbf{C}_k &= \mathbf{D}_{k+1} [\mathbf{P}_{k+1}^-]^{-1} \\
\mathbf{m}_k^s &= \mathbf{m}_k + \mathbf{C}_k (\mathbf{m}_{k+1}^s - \mathbf{m}_{k+1}^-) \\
\mathbf{P}_k^s &= \mathbf{P}_k + \mathbf{C}_k (\mathbf{P}_{k+1}^s - \mathbf{P}_{k+1}^-) \mathbf{C}_k^T.
\end{aligned} \tag{4.32}$$

The integrals above can be approximated using analogous numerical integration or analytical approximation schemes as in the filtering case, that is, with Gauss-Hermite quadratures or central differences (Ito and Xiong, 2000; Nørgaard et al., 2000; Wu et al., 2006), cubature rules (Arasaratnam and Haykin, 2009), Monte Carlo (Kotecha and Djuric, 2003), Gaussian process / Bayes-Hermite based integration (O'Hagan, 1991; Deisenroth et al., 2009), or with many other numerical integration schemes.

4.3 Monte Carlo Based Smoothers

4.3.1 Sequential Importance Resampling Smoother

Optimal smoothing can be performed with the SIR algorithm with a slight modification to the filtering case. Instead of keeping Monte Carlo samples of the states on single time step $\mathbf{x}_k^{(i)}$, we keep samples of the whole state histories $\mathbf{x}_{1:k}^{(i)}$. The computations of the algorithm remain exactly the same, but in resampling stage the whole state histories are resampled instead of the states of single time steps. The weights of these state histories are the same as in normal SIR algorithm and the smoothed posterior distribution estimate of time step k given the measurements up to the time step $T > k$ is given as (Kitagawa, 1996; Doucet et al., 2000)

$$p(\mathbf{x}_k | \mathbf{y}_{1:T}) \approx \sum_{i=1}^N w_T^{(i)} \delta(\mathbf{x}_k - \mathbf{x}_k^{(i)}). \tag{4.33}$$

where $\delta(\cdot)$ is the Dirac delta function and $\mathbf{x}_k^{(i)}$ is the k th component in $\mathbf{x}_{1:T}^{(i)}$.

However, if $T \gg k$ this simple method is known to produce very degenerate approximations (Kitagawa, 1996; Doucet et al., 2000). In (Godsill et al., 2004) more efficient methods for sampling from the smoothing distributions are presented.

4.3.2 Rao-Blackwellized Particle Smoother

The Rao-Blackwellized particle smoother can be used for computing the smoothing solution to the conditionally Gaussian RBPF model (3.124). A weighted set of Monte Carlo samples from the smoothed distribution of the parameters θ_k in the model (3.124) can be produced by storing the histories instead of the single states, as in the case of plain SIR. The corresponding histories of the means and the covariances are then conditional on the *parameter histories* $\theta_{1:T}$. However, the means and covariances at time step k are only conditional on the *measurement histories* up to k , not on the later measurements. In order to correct this, Kalman smoothers have to be applied to each history of the means and the covariances.

Algorithm 4.6 (Rao-Blackwellized particle smoother). *A set of weighted samples $\{w_T^{s,(i)}, \theta_{1:T}^{s,(i)}, \mathbf{m}_{1:T}^{s,(i)}, \mathbf{P}_{1:T}^{s,(i)} : i = 1, \dots, N\}$ representing the smoothed distribution can be computed as follows:*

1. Compute the weighted set of Rao-Blackwellized state histories

$$\{w_T^{(i)}, \theta_{1:T}^{(i)}, \mathbf{m}_{1:T}^{(i)}, \mathbf{P}_{1:T}^{(i)} : i = 1, \dots, N\} \quad (4.34)$$

by using the Rao-Blackwellized particle filter.

2. Set

$$\begin{aligned} w_T^{s,(i)} &= w_T^{(i)} \\ \theta_{1:T}^{s,(i)} &= \theta_{1:T}^{(i)}. \end{aligned} \quad (4.35)$$

3. Apply the Kalman smoother to each of the mean and covariance histories $\mathbf{m}_{1:T}^{(i)}, \mathbf{P}_{1:T}^{(i)}$ for $i = 1, \dots, N$ to produce the smoothed mean and covariance histories $\mathbf{m}_{1:T}^{s,(i)}, \mathbf{P}_{1:T}^{s,(i)}$.

The Rao-Blackwellized particle smoother in this simple form also has the same disadvantage as the plain SIR smoother, that is, the smoothed estimate of θ_k can be quite degenerate if $T \gg k$. Fortunately, the smoothed estimates of the actual states \mathbf{x}_k can still be quite good, because its degeneracy is avoided by the Rao-Blackwellization. To avoid the degeneracy in estimates of θ_k it is possible to use more efficient sampling procedures for generating samples from the smoothing distributions (Fong et al., 2002).

As in the case of filtering, in some cases approximately Gaussian parts of a state space model can be approximately marginalized by using extended Kalman smoothers or unscented Kalman smoothers.

In the case of Rao-Blackwellization of static parameters (Storvik, 2002) the smoothing is much easier. In this case, due to lack of dynamics, the posterior distribution obtained after processing the last measurement is the smoothed distribution.

Appendix A

Additional Material

A.1 Properties of Gaussian Distribution

Definition A.1 (Gaussian distribution). *Random variable $\mathbf{x} \in \mathbb{R}^n$ has Gaussian distribution with mean $\mathbf{m} \in \mathbb{R}^n$ and covariance $\mathbf{P} \in \mathbb{R}^{n \times n}$ if it has the probability density of the form*

$$N(\mathbf{x} | \mathbf{m}, \mathbf{P}) = \frac{1}{(2\pi)^{n/2} |\mathbf{P}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \mathbf{P}^{-1} (\mathbf{x} - \mathbf{m})\right), \quad (\text{A.1})$$

where $|\mathbf{P}|$ is the determinant of matrix \mathbf{P} .

Lemma A.1 (Joint density of Gaussian variables). *If random variables $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in \mathbb{R}^m$ have the Gaussian probability densities*

$$\begin{aligned} \mathbf{x} &\sim N(\mathbf{x} | \mathbf{m}, \mathbf{P}) \\ \mathbf{y} | \mathbf{x} &\sim N(\mathbf{y} | \mathbf{H}\mathbf{x} + \mathbf{u}, \mathbf{R}), \end{aligned} \quad (\text{A.2})$$

then the joint density of \mathbf{x}, \mathbf{y} and the marginal distribution of \mathbf{y} are given as

$$\begin{aligned} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} &\sim N\left(\begin{bmatrix} \mathbf{m} \\ \mathbf{H}\mathbf{m} + \mathbf{u} \end{bmatrix}, \begin{bmatrix} \mathbf{P} & \mathbf{P}\mathbf{H}^T \\ \mathbf{H}\mathbf{P} & \mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R} \end{bmatrix}\right) \\ \mathbf{y} &\sim N(\mathbf{H}\mathbf{m} + \mathbf{u}, \mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R}). \end{aligned} \quad (\text{A.3})$$

Lemma A.2 (Conditional density of Gaussian variables). *If the random variables \mathbf{x} and \mathbf{y} have the joint Gaussian probability density*

$$\mathbf{x}, \mathbf{y} \sim N\left(\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix}\right), \quad (\text{A.4})$$

then the marginal and conditional densities of \mathbf{x} and \mathbf{y} are given as follows:

$$\begin{aligned} \mathbf{x} &\sim N(\mathbf{a}, \mathbf{A}) \\ \mathbf{y} &\sim N(\mathbf{b}, \mathbf{B}) \\ \mathbf{x} | \mathbf{y} &\sim N(\mathbf{a} + \mathbf{C}\mathbf{B}^{-1}(\mathbf{y} - \mathbf{b}), \mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^T) \\ \mathbf{y} | \mathbf{x} &\sim N(\mathbf{b} + \mathbf{C}^T \mathbf{A}^{-1}(\mathbf{x} - \mathbf{a}), \mathbf{B} - \mathbf{C}^T \mathbf{A}^{-1} \mathbf{C}). \end{aligned} \quad (\text{A.5})$$

References

- Akashi, H. and Kumamoto, H. (1977). Random sampling approach to state estimation in switching environments. *Automatica*, 13:429–434.
- Anderson, R. M. and May, R. M. (1991). *Infectious Diseases of Humans: Dynamics and Control*. Oxford University Press, Oxford.
- Andrieu, C., de Freitas, N., and Doucet, A. (2002). Rao-Blackwellised particle filtering via data augmentation. In Dietterich, T. G., Becker, S., and Ghahramani, Z., editors, *Advances in Neural Information Processing Systems 14*. MIT Press.
- Arasaratnam, I. and Haykin, S. (2009). Cubature Kalman filters. *IEEE Transactions on Automatic Control*, 54(6):1254–1269.
- Bar-Shalom, Y. and Li, X.-R. (1995). *Multitarget-Multisensor Tracking: Principles and Techniques*. YBS.
- Bar-Shalom, Y., Li, X.-R., and Kirubarajan, T. (2001). *Estimation with Applications to Tracking and Navigation*. Wiley, New York.
- Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer.
- Bernardo, J. M. and Smith, A. F. M. (1994). *Bayesian Theory*. John Wiley & Sons.
- Bishop, C. M. (1995). *Neural Networks for Pattern Recognition*. Oxford University Press.
- Blackman, S. and Popoli, R. (1999). *Design and Analysis of Modern Tracking Systems*. Artech House Radar Library.
- Casella, G. and Robert, C. P. (1996). Rao-Blackwellisation of sampling schemes. *Biometrika*, 83(1):81–94.
- Cox, H. (1964). On the estimation of state variables and parameters for noisy dynamic systems. *IEEE Transactions on Automatic Control*, 9(1):5–12.
- Deisenroth, M. P., Huber, M. F., and Hanebeck, U. D. (2009). Analytic moment-based Gaussian process filtering. In *Proceedings of the 26th International Conference on Machine Learning*.
- Doucet, A., de Freitas, N., and Gordon, N. (2001). *Sequential Monte Carlo Methods in Practice*. Springer, New York.
- Doucet, A., Godsill, S. J., and Andrieu, C. (2000). On sequential Monte Carlo sampling methods for Bayesian filtering. *Statistics and Computing*, 10(3):197–208.
- Fearnhead, P. and Clifford, P. (2003). On-line inference for hidden Markov models via particle filters. *J. R. Statist. Soc. B*, 65(4):887–899.
- Fong, W., Godsill, S. J., Doucet, A., and West, M. (2002). Monte Carlo smoothing with application to audio signal enhancement. *IEEE Transactions on Signal Processing*, 50(2):438–449.
- Gelb, A. (1974). *Applied Optimal Estimation*. The MIT Press, Cambridge, MA.
- Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. R. (1995). *Bayesian Data Analysis*. Chapman & Hall.
- Gilks, W., Richardson, S., and Spiegelhalter, D., editors (1996). *Markov Chain Monte*

- Carlo in Practice*. Chapman & Hall.
- Godsill, S. J., Doucet, A., and West, M. (2004). Monte Carlo smoothing for nonlinear time series. *Journal of the American Statistical Association*, 99(465):156–168.
- Godsill, S. J. and Rayner, P. J. (1998). *Digital Audio Restoration: A Statistical Model Based Approach*. Springer-Verlag.
- Gordon, N. J., Salmond, D. J., and Smith, A. F. M. (1993). Novel approach to nonlinear/non-Gaussian Bayesian state estimation. In *IEEE Proceedings on Radar and Signal Processing*, volume 140, pages 107–113.
- Grewal, M. S. and Andrews, A. P. (2001). *Kalman Filtering, Theory and Practice Using MATLAB*. Wiley, New York.
- Grewal, M. S., Weill, L. R., and Andrews, A. P. (2001). *Global Positioning Systems, Inertial Navigation and Integration*. Wiley, New York.
- Hayes, M. H. (1996). *Statistical Digital Signal Processing and Modeling*. John Wiley & Sons, Inc.
- Ho, Y. C. and Lee, R. C. K. (1964). A Bayesian approach to problems in stochastic estimation and control. *IEEE Transactions on Automatic Control*, 9:333–339.
- Ito, K. and Xiong, K. (2000). Gaussian filters for nonlinear filtering problems. *IEEE Transactions on Automatic Control*, 45(5):910–927.
- Jazwinski, A. H. (1966). Filtering for nonlinear dynamical systems. *IEEE Transactions on Automatic Control*, 11(4):765–766.
- Jazwinski, A. H. (1970). *Stochastic Processes and Filtering Theory*. Academic Press, New York.
- Julier, S. J. and Uhlmann, J. K. (1995). A general method of approximating nonlinear transformations of probability distributions. Technical report, Robotics Research Group, Department of Engineering Science, University of Oxford.
- Julier, S. J. and Uhlmann, J. K. (2004). Unscented filtering and nonlinear estimation. *Proceedings of the IEEE*, 92(3):401–422.
- Julier, S. J., Uhlmann, J. K., and Durrant-Whyte, H. F. (1995). A new approach for filtering nonlinear systems. In *Proceedings of the 1995 American Control Conference, Seattle, Washington*, pages 1628–1632.
- Julier, S. J., Uhlmann, J. K., and Durrant-Whyte, H. F. (2000). A new method for the nonlinear transformation of means and covariances in filters and estimators. *IEEE Transactions on Automatic Control*, 45(3):477–482.
- Kaipio, J. and Somersalo, E. (2005). *Statistical and Computational Inverse Problems*. Number 160 in Applied mathematical Sciences. Springer.
- Kalman, R. E. (1960a). Contributions to the theory of optimal control. *Boletin de la Sociedad Matematica Mexicana*, 5(1):102–119.
- Kalman, R. E. (1960b). A new approach to linear filtering and prediction problems. *Transactions of the ASME, Journal of Basic Engineering*, 82:35–45.
- Kalman, R. E. and Bucy, R. S. (1961). New results in linear filtering and prediction theory. *Transactions of the ASME, Journal of Basic Engineering*, 83:95–108.
- Kaplan, E. D. (1996). *Understanding GPS, Principles and Applications*. Artech House, Boston, London.
- Kitagawa, G. (1996). Monte Carlo filter and smoother for non-Gaussian nonlinear state space models. *Journal of Computational and Graphical Statistics*, 5:1–25.
- Kotecha, J. H. and Djuric, P. M. (2003). Gaussian particle filtering. *IEEE Transactions on Signal Processing*, 51(10).
- Lee, R. C. K. (1964). *Optimal Estimation, Identification and Control*. M.I.T. Press.
- Lefebvre, T., Bruyninckx, H., and Schuller, J. D. (2002). Comment on "a new method for

- the nonlinear transformation of means and covariances in filters and estimators" [and authors' reply]. *IEEE Transactions on Automatic Control*, 47(8):1406–1409.
- Liu, J. S. (2001). *Monte Carlo Strategies in Scientific Computing*. Springer.
- Liu, J. S. and Chen, R. (1995). Blind deconvolution via sequential imputations. *Journal of the American Statistical Association*, 90(430):567–576.
- MacKay, D. J. C. (1998). Introduction to Gaussian processes. In Bishop, C. M., editor, *Neural Networks and Machine Learning*, pages 133–165. Springer-Verlag.
- Maybeck, P. (1982a). *Stochastic Models, Estimation and Control, Volume 2*. Academic Press.
- Maybeck, P. (1982b). *Stochastic Models, Estimation and Control, Volume 3*. Academic Press.
- Milton, J. S. and Arnold, J. C. (1995). *Introduction to Probability and Statistics, Principles and Applications for Engineering and the Computing Sciences*. McGraw-Hill, Inc.
- Murray, J. D. (1993). *Mathematical Biology*. Springer, New York.
- Nørgaard, M., Poulsen, N. K., and Ravn, O. (2000). New developments in state estimation for nonlinear systems. *Automatica*, 36(11):1627 – 1638.
- O'Hagan, A. (1991). Bayes-Hermite quadrature. *Journal of Statistical Planning and Inference*, 29:245–260.
- Pikkarainen, H. (2005). *A Mathematical Model for Electrical Impedance Process Tomography*. Doctoral dissertation, Helsinki University of Technology.
- Pitt, M. K. and Shephard, N. (1999). Filtering via simulation: Auxiliary particle filters. *Journal of the American Statistical Association*, 94(446):590–599.
- Proakis, J. G. (2001). *Digital Communications*. McGraw-Hill, 4th edition.
- Punskaya, E., Doucet, A., and Fitzgerald, W. J. (2002). On the use and misuse of particle filtering in digital communications. In *EUSIPCO*.
- Rafael C. Gonzalez, R. E. W. (2008). *Digital Image Processing*. Prentice Hall, 3rd edition.
- Raiffa, H. and Schlaifer, R. (2000). *Applied Statistical Decision Theory*. John Wiley & Sons, Wiley Classics Library.
- Rauch, H. E., Tung, F., and Striebel, C. T. (1965). Maximum likelihood estimates of linear dynamic systems. *AIAA Journal*, 3(8):1445–1450.
- Ristic, B., Arulampalam, S., and Gordon, N. (2004). *Beyond the Kalman Filter*. Artech House, Norwood, MA.
- Sage, A. P. and Melsa, J. L. (1971). *Estimation Theory with Applications to Communications and Control*. McGraw-Hill Book Company.
- Särkkä, S. (2008). Unscented Rauch-Tung-Striebel smoother. *IEEE Transactions on Automatic Control*, 53(3):845–849.
- Särkkä, S. and Hartikainen, J. (2010a). On Gaussian optimal smoothing of non-linear state space models. (submitted for publication).
- Särkkä, S. and Hartikainen, J. (2010b). Sigma point methods in optimal smoothing of non-linear stochastic state space models. (submitted for publication).
- Shiryayev, A. N. (1996). *Probability*. Springer.
- Stengel, R. F. (1994). *Optimal Control and Estimation*. Dover Publications, New York.
- Stone, L. D., Barlow, C. A., and Corwin, T. L. (1999). *Bayesian Multiple Target Tracking*. Artech House, Boston, London.
- Storvik, G. (2002). Particle filters in state space models with the presence of unknown static parameters. *IEEE Transactions on Signal Processing*, 50(2):281–289.
- Stratonovich, R. L. (1968). *Conditional Markov Processes and Their Application to the Theory of Optimal Control*. American Elsevier Publishing Company, Inc.
- Titterton, D. H. and Weston, J. L. (1997). *Strapdown Inertial Navigation Technology*. Peter

- Pregrinus Ltd.
- Van der Merwe, R., Freitas, N. D., Doucet, A., and Wan, E. (2001). The unscented particle filter. In *Advances in Neural Information Processing Systems 13*.
- van der Merwe, R. and Wan, E. (2003). Sigma-point Kalman filters for probabilistic inference in dynamic state-space models. In *Proceedings of the Workshop on Advances in Machine Learning*.
- Van Trees, H. L. (1968). *Detection, Estimation, and Modulation Theory Part I*. John Wiley & Sons, New York.
- Van Trees, H. L. (1971). *Detection, Estimation, and Modulation Theory Part II*. John Wiley & Sons, New York.
- Vauhkonen, M. (1997). *Electrical impedance tomography and prior information*. PhD thesis, Kuopio University.
- Viterbi, A. J. (1967). Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory*, II-13(2).
- Wan, E. A. and Van der Merwe, R. (2001). The unscented Kalman filter. In Haykin, S., editor, *Kalman Filtering and Neural Networks*, chapter 7. Wiley.
- West, M. and Harrison, J. (1997). *Bayesian Forecasting and Dynamic Models*. Springer-Verlag.
- Wiener, N. (1950). *Extrapolation, Interpolation and Smoothing of Stationary Time Series with Engineering Applications*. John Wiley & Sons, Inc., New York.
- Wu, Y., Hu, D., Wu, M., and Hu, X. (2005). Unscented Kalman filtering for additive noise case: Augmented versus nonaugmented. *IEEE Signal Processing Letters*, 12(5):357–360.
- Wu, Y., Hu, D., Wu, M., and Hu, X. (2006). A numerical-integration perspective on Gaussian filters. *IEEE Transactions on Signal Processing*, 54(8):2910–2921.