MATH4406 (Control Theory) Part 10: Kalman Filtering and LQG Prepared by Yoni Nazarathy, Last Updated: October 19, 2012

1 About

We have spent plenty of time in the course dealing with systems of the form:

$$\dot{x}(t) = Ax(t) + Bu(t) \\ y(t) = Cx(t) + Du(t)$$
 and
$$\begin{aligned} x(n+1) &= Ax(n) + Bu(n) \\ y(n) &= Cx(n) + Du(n) \end{aligned}$$
(1)

with $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times m}$, $C \in \mathbb{R}^{p \times n}$ and $D \in \mathbb{R}^{p \times m}$. The focus was mostly on the continuous time version (u(t), x(t), y(t)). In unit 4 we saw how to design a state feedback controller and an observer and in later units we dealt with optimal control of such systems.

We now augment our system models by adding *noise* components. To the first equation we shall add *disturbance noise* (ξ_x) and to the second equation we shall add *measurement noise* (ξ_y) . This yields:

$$\dot{x}(t) = Ax(t) + Bu(t) + \xi_x(t)$$

$$y(t) = Cx(t) + Du(t) + \xi_y(t)$$
 or
$$x(n+1) = Ax(n) + Bu(n) + \xi_x(n)$$

$$y(n) = Cx(n) + Du(n) + \xi_y(n)$$

One way of modeling the noise is by assuming that $\xi(\cdot)$ is from some function class and assuming that in controlling the system we have no knowledge of what specific $\xi(\cdot)$ from that class is used. This is the method of *robust control*. Alternatively, we can think of $\xi(\cdot)$ as a *random process(es)* by associating a probability space with the model. We shall focus on the latter approach.

The level of mathematical care that is needed for handling the continuous time case is beyond our scope as it requires some basic understanding of *stochastic calculus* (e.g. Brownian motion, Stochastic Differential Equations, Ito's formula etc...). We shall thus focus on the discrete time case which is simpler in that the random processes (es) are discrete time sequences of *random variables*. Luckily, the methods that we shall survey (Kalman filtering and Linear Quadratic Gaussian (LQG) optimal control) are often applied in practice in discrete time on digital computers. So understanding the discrete time case is both pedagogically simpler and often of greater practical interest.

In treating $\xi_x(n)$ and $\xi_y(n)$ as discrete time random processes we shall assume they are each i.i.d. (independent and identically distributed) sequences of zero-mean Gaussian (normal) random vectors with covariance matrices Σ_x and Σ_y respectively (we review this below). In many physical situations this is actually a practical assumption:

- Having the noise of one time slot independent of the disturbance at other time slots is the practical situation (especially for short lived disturbances). (This is the first 'i' of i.i.d.).
- Having noise of a constant statistical law makes sense for time invariant systems. (This is the second 'i' of i.i.d.).
- Having noise that have a mean of 0 implies there is no general direction of the disturbance.
- Having noise that follows the Gaussian distribution is sensible if the noise is a summation of many small factors. In this case the *central limit theorem* implies that the noise distribution is Gaussian.

Note 1: We are not restricting individual coordinates of $\xi(n)$ (at any time n) to be independent.

Note 2: Note that even though the noise terms are i.i.d., $x(\cdot)$ is no longer an i.i.d. process (it will be in the pathological case in which A = 0 and B = 0).

Note 3: In many situations the variance (covariance matrix) of ξ can be modeled from "first principles" just as the (A, B, C, D) model is. This is the case of noise is due to well understood electromagnetic effects as well as due to rounding errors appearing in digital control.

What will we do with the stochastic model?

1. State estimation (Kalman filtering): For the deterministic system, we saw the Luenberger observer:

$$\hat{x}(n+1) = A\hat{x}(n) + Bu(n) + K(y(n) - \hat{y}(n)).$$

The Kalman filter is used to do essentially the same thing, yet now taking into control the fact that now $x(\cdot)$ is a random process.

2. Optimal control (LQG): For the deterministic system we saw how to design a state feedback control such that,

$$\sum_{k=0}^{\infty} x'(k)Qx(k) + u'(k)Ru(k),$$

is minimized (if $Q \ge 0$ and R > 0). Now with random noise, $x(\cdot)$ is a random process. Further if we use a state feedback control then $u(\cdot)$ is random process. We are thus interested in finding a control law that minimizes,

$$\mathbb{E}\left[\sum_{k=0}^{\infty} x'(k)Qx(k) + u'(k)Ru(k)\right]$$

We will have time to touch LQG only briefly and focus mostly on the Kalman filter. A practical note: The celebrated Kalman filter is implemented in a variety of engineering applications dealing with tracking, positioning and sensing. It is a good thing to know about outside the scope of control also.

2 Gaussian Random Vectors

We briefly review/summarize Gaussian random vectors. We begin with Gaussian scalars: A random variable, X is said to have a Gaussian (normal) distribution with a mean of μ and a variance of $\sigma^2 > 0$, denoted, $X \sim N(\mu, \sigma^2)$ if,

$$\mathbb{P}\left(a \le X \le b\right) = \int_{a}^{b} \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^{2}} dx.$$

We have,

$$\mathbb{E}\left[X\right] = \int_{-\infty}^{\infty} x \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx = \mu.$$

Further,

$$Var(X) = \int_{-\infty}^{\infty} (x-\mu)^2 \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} dx = \sigma^2.$$

We now consider the random vector $X = [X_1, \ldots, X_n]'$. Assume momentarily that each of the random variables X_i follows some arbitrary distribution. The expectation of X is the vector,

$$\mathbb{E}\left[X\right] = \left[\begin{array}{c} \mathbb{E}\left[X_1\right]\\ \vdots\\ \mathbb{E}\left[X_n\right] \end{array}\right].$$

Similarly, if X is a matrix of random variables than $\mathbb{E}[X]$ is the matrix of the expectations of the individual entries.

This leads to the important (for our context) definition of the *covariance matrix* of a random vector X.

$$Cov(X) = \mathbb{E} \left[\left(X - \mathbb{E} \left[X \right] \right) \left(X - \mathbb{E} \left[X \right] \right)' \right].$$

Note that Cov(X) is an $n \times n$ symmetric matrix with individual elements:

$$\left(Cov(X)\right)_{i,j} = Cov(X_i, X_j).$$

Reminder: For two scalar random variables Z and W,

$$Cov(Z,W) = \mathbb{E}\left[\left(Z - \mathbb{E}\left[Z\right]\right)\left(W - \mathbb{E}\left[W\right]\right)\right] = \mathbb{E}\left[ZW\right] - \mathbb{E}\left[Z\right]\mathbb{E}\left[W\right].$$

Notes: (i) If Z = W then Cov(Z, W) = Var(Z). (ii) If one (or both) of the random variables are zero mean then the covariance is simply $\mathbb{E}[ZW]$. (iii) If the random variables are independent then the covariance is 0.

The covariance matrix thus records the variance of the random variables on the diagonal and the covariances on the off-diagonal entries.

Exercise 2.1 Assume you are given an n dimensional random variable X and an $m \times n$ matrix A. Define Y = AX. What is the mean vector of Y in terms of that of X? What is the covariance matrix of Y in terms of the covariance matrix of X?

We are now in a position to define (one of several equivalent definitions) Gaussian random vectors: We say the random vector X is Gaussian with mean vector μ and covariance matrix Σ , denoted $X \sim N(\mu, \Sigma)$ if,

$$\mathbb{P}\Big(a_1 \le X_1 \le b_1, \dots, a_n \le X_n \le b_n\Big) = \int_{a_1}^{b_1} \dots \int_{a_n}^{b_n} \phi(x_1, \dots, x_n) dx_1 \dots dx_n,$$

with the density function being,

$$\phi(x_1,\ldots,x_n) = \frac{1}{(2\pi)^{n/2} det(\Sigma)^{1/2}} e^{-\frac{1}{2}(x-\mu)'\Sigma^{-1}(x-\mu)}$$

It can be calculated that in this case,

$$\mathbb{E}[X] = \mu, \qquad Cov(X) = \Sigma.$$

Further, the marginal distribution of each of the X_i 's is normal.

Notes: (i) Distributions of Gaussian random vectors are characterized by their mean vector and covariance matrix. (ii) If two coordinates are non-correlated (covariance 0) then they are independent. (iii) Linear transformations of Gaussian random vectors yield random vectors that still follow the Gaussian distribution with mean and covariance as given by Exercise 2.1.

The final property that we shall overview for Gaussian random vectors deals with conditional distributions. Partition $X \sim \mathcal{N}(\mu, \Sigma)$ into X_a and X_b and have,

$$\mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix} . \qquad \Sigma = \begin{bmatrix} \Sigma_a & \Sigma_{ab} \\ \Sigma'_{ab} & \Sigma_b \end{bmatrix}$$

We have that the distribution of X_a conditional on $X_b = x_b$ is

$$\mathcal{N}\Big(\mu_a + \Sigma_{ab}\Sigma_b^{-1}(x_b - \mu_b), \quad \Sigma_a - \Sigma_{ab}\Sigma_b^{-1}\Sigma_{ab}'\Big). \tag{2}$$

This is useful for estimating X_a based on *measurements* of X_b . A sensible estimate in this case is, $\mu_a + \sum_{ab} \sum_{b}^{-1} (x_b - \mu_b)$. As a "sanity check" of this formula observe that if X_a and X_b are independent then $\sum_{ab} = 0$ and thus the estimate is simply μ_a .

3 Minimum Mean Square Estimation

Consider now the general situation in which you observe the value of a random vector $X_b = x_b$ and would like to use it to estimate the value of X_a . Here we model X_a and X_b as two random vectors (measurable functions) on the same probability space and hope that they are somewhat dependent (i.e. knowledge of X_b can give us some information on X_a). We are thus looking for a function $f(\cdot)$ such that $f(X_b)$ is a "good" estimate on X_a . There are all kinds of definitions of "good" – here is perhaps the most popular one:

$$\min_{h} \mathbb{E}\left[||X_a - h(X_b)||^2 \right], \tag{3}$$

where $||\cdot||$ is the Euclidean norm and the minimization is over all $h(\cdot)$ in some function class whose definition we leave vague for the purpose of this informal discussion. Note that the expectation is with respect to both X_a and X_b . Does this criterion make sense? Yes, of course! Further, it turns out to be very tractable in certain cases since it turns out that the $h(\cdot)$ that minimizes (3) is:

$$h^*(x_b) = \mathbb{E} \left[X_a \mid X_b = x_b \right]. \tag{4}$$

The above is read as the "conditional expectation of the random vector X_a , given the observed value x_b ". Does the best estimator $h^*(\cdot)$ make sense? Yes of course!

Brief reminder: If two random vectors X_a and X_b are distributed say with a density $f_{ab}(x_a, x_b)$, then the conditional density of X_a given $X_b = x_b$ is:

$$f_{a|b}(x_a|x_b) = \frac{f_{ab}(x_a, x_b)}{f_b(x_b)}$$

where the denominator is the marginal density of X_b , namely (assuming X_a is k-dimensional):

$$f_b(x_b) = \int_{x_a \in \mathbb{R}^k} f_{ab}(x_a, x_b) dx_a.$$

I.e. to get the marginal density of X_b you need to "integrate out" all of the values that X_a may get. And to get the conditional distribution of X_a given the information that X_b takes a specific values x_b , you need to "rescale" the joint density by the marginal of X_b . Try to draw this in two dimensions.

Now the conditional expectation (for a given value of X_b) that appears in (4) is simply evaluated as follows:

$$\mathbb{E}\left[X_a \mid X_b = x_b\right] = \int_{x_a \in \mathbb{R}^k} x_a \ f_{a|b}(x_a|x_b) dx_a.$$

Further note that the expression $\mathbb{E} [X_a | X_b]$ (where we do not specify a given values for X_b) is actually a random variable that is a function of the random variable X_b , where the function is:

$$g(X_b) = \int_{x_a \in \mathbb{R}^k} x_a \ f_{a|b}(x_a|X_b) dx_a.$$

Hence the conditional expectation $\mathbb{E} [X_a \mid X_b]$ is actually a random variable in itself. And we may thus attempt to take its expectation. It turns out that in this case:

$$\mathbb{E}\left[g(X_b)\right] = \mathbb{E}\left[\mathbb{E}\left[X_a \mid X_b\right]\right] = \mathbb{E}\left[X_a\right].$$
(5)

Note: The above "brief reminder" about conditional expectation is very informal as technical details are missing. Yet this is enough for our needs.

Here is now (an informal) proof of (4):

Proof First use the conditional expectation formula similar to (5):

$$\mathbb{E}\left[||X_a - h(X_b)||^2\right] = \mathbb{E}\left[\mathbb{E}\left[||X_a - h(X_b)||^2 \mid X_b\right]\right] = \int \mathbb{E}\left[||X_a - h(X_b)||^2 \mid X_b(\omega)\right] dP_{X_b}(\omega)$$
(6)

The last expression represents the outer expectation as a Lebesgue integral with respect to the probability measure associated with the random variable X_b . This is not needed to understand the proof, but is here for additional clarity on the meaning of expectation.

Note that the internal conditional expectation (conditional on X_b) is a function, $\tilde{g}(\cdot)$ of the random variable X_b . Let's investigate this function in a bit greater detail. Assume that the estimator $h(X_b)$ takes on the value z (i.e. assume that in the probability sample space associated with the random variable X_b , we get and ω such that $h(X_b(\omega)) = z$). Then,

$$\mathbb{E}\left[||X_a - z||^2 \mid X_b\right] = \mathbb{E}\left[||X_a||^2 \mid X_b\right] - 2z'\mathbb{E}\left[X_a \mid X_b\right] + ||z||^2$$

Taking derivative with respect to z (note that z is generally a vector) and equating to 0 implies that the above is minimized by $z = \mathbb{E} [X_a | X_b]$. I.e. the integrand in (6) is minimized by setting,

$$h(X_b) = \mathbb{E}\left[X_a | X_b\right].$$

Thus the integral (the outer expectation) is also minimized by this choice of $h(\cdot)$ and thus the (3) is minimized by (4).

Evaluating (4) for arbitrarily distributed X_a and X_b can be a complicated (not explicitly solvable) task. Yet for Gaussian random vectors we are blessed with a clean result. Indeed as we saw in the case of Gaussian random vectors that this conditional expectation has the closed (linear) form. So if you believe (4), in the case of Gaussian random vectors,

$$h^*(x_b) = \mu_a + \Sigma_{ab} \Sigma_b^{-1}(x_b - \mu_b).$$

We thus see that for Gaussian random vectors, the optimal estimator $h^*(\cdot)$ is an linear (affine to be precise) function of x_b . It is thus tempting to restrict the function class of $h^*(\cdot)$ in (3) to,

$$h(x_b) = Gx_b + g,$$

where G and g are a matrix and a vector of the appropriate dimension. The pair (G, g) that minimizes (3) is sometimes called the LMMSE estimator (Linear Minimum Mean Square Error estimator).

Exercise 3.1 What are G and g in the case of Gaussian random variables?

Exercise 3.2 Prove the following proposition by taking derivatives w.r.t. to G and g.

Proposition 3.3 Let (X_a, X_b) be random vectors with means μ_a and μ_b respectively and with a covariance matrix (of $(X_a, X_b)'$) being:

$$\begin{bmatrix} \Sigma_a & \Sigma_{ab} \\ \Sigma'_{ab} & \Sigma_b \end{bmatrix}$$

Then LMMSE estimator of X_a given $X_b = x_b$ is:

$$h^*(x_b) = \mu_a + \Sigma_{ab} \Sigma_b^{-1}(x_b - \mu_b).$$

Further the covariance matrix of the error vector $X_a - h^*(X_b)$ is given by:

$$\mathbb{E}\left[\left(X_a - h^*(X_b)\right)\left(X_a - h^*(X_b)\right)'\right] = \Sigma_a - \Sigma_{a,b}\Sigma_b^{-1}\Sigma'_{a,b}$$

In the case of non-Gaussian random variables, restricting to an affine estimator based on G and g is often a compromise:

Exercise 3.4 Let X_b have a uniform distribution on the interval [-1, 1] and let $X_a = X_b^2$. Find the best affine estimator of X_a in terms of X_b and compare its performance (using the objective (3)) to the best estimator (4).

Repeat for the case of,

$$f_{a,b}(x_a, x_b) = \begin{cases} 2e^{-(x_a + x_b)} & 0 \le x_b \le x_a < \infty, \\ 0 & elsewhere. \end{cases}$$

4 The Kalman Filtering Problem "Solved" by LMMSE

Our goal is to have a state estimate, $\hat{x}(n)$ for a given (A, B, C, D) + noise system:

$$\begin{aligned} x(n+1) &= Ax(n) + Bu(n) + \xi_x(n) \\ y(n) &= Cx(n) + Du(n) + \xi_y(n) \end{aligned}$$

More specifically we assume we have controlled this system over times k = 0, ..., N - 1 by setting inputs u(0), ..., u(N-1) (which we know) and have measured outputs y(0), ..., y(N-1). Note that we treat x(0) as a random variable also where we assume we know its mean and covariance.

We will now show that his problem can be posed as estimating X_a based on measurement of X_b (as presented in the previous section) where,

$$X_a = (x(0)', x(1)', \dots, x(N)')', \qquad X_b = (y(0)', y(1)', \dots, y(N)')',$$

and the inputs $u(0), \ldots, u(N-1)$ are known values.

By iterating the system, we get:

$$\begin{aligned} x(1) &= Ax(0) + Bu(0) + \xi_x(0), \\ x(2) &= A^2 x(0) + ABu(0) + Bu(1) + A\xi_x(0) + \xi_x(1), \\ x(3) &= A^3 x(0) + A^2 Bu(0) + ABu(1) + Bu(2) + A^2 \xi_x(0) + A\xi_x(1) + \xi_x(2), \\ &\vdots \\ x(N) &= A^N x(0) + \sum_{k=0}^{N-1} A^{N-1-k} Bu(k) + \sum_{k=0}^{N-1} A^{N-1-k} \xi_x(k). \end{aligned}$$

Plugging the above in the output equations, we get,

$$\begin{array}{lll} y(0) &=& Cx(0) + Du(0) + \xi_y(0), \\ y(1) &=& CAx(0) + CBu(0) + C\xi_x(0) + Du(1) + \xi_y(1) \\ y(2) &=& CA^2x(0) + CABu(0) + CBu(1) + CA\xi_x(0) + C\xi_x(1) + Du(2) + \xi_y(2) \\ &\vdots \\ y(N) &=& CA^Nx(0) + \sum_{k=0}^{N-1} (CA^{N-1-k}B)u(k) + Du(N) + \sum_{k=0}^{N-1} CA^{N-1-k}\xi_x(k) + \xi_y(N) \end{array}$$

It is thus a simple matter to write out constant matrices \tilde{A}, \tilde{C} and well as functions of

the known input, $\tilde{b}(u), \tilde{d}(u)$, such that:

$$X_{a} = \begin{bmatrix} x(0) \\ x(1) \\ \vdots \\ x(N) \end{bmatrix} = \tilde{A} \begin{bmatrix} x(0) \\ \xi_{x}(0) \\ \vdots \\ \xi_{x}(N-1) \end{bmatrix} + \tilde{b}(u(0), \dots, u(N-1)),$$
$$X_{b} = \begin{bmatrix} y(0) \\ \vdots \\ y(N) \end{bmatrix} = \tilde{C} \begin{bmatrix} x(0) \\ \xi_{x}(0) \\ \vdots \\ \xi_{x}(N-1) \end{bmatrix} + \begin{bmatrix} \xi_{y}(0) \\ \vdots \\ \xi_{y}(N) \end{bmatrix} + \tilde{d}(u(0), \dots, u(N))$$

Exercise 4.1 Specify \tilde{A}, \tilde{C} as well as $\tilde{b}(u), \tilde{d}(u)$ explicitly.

It is now useful to consider the combined random vector,

$$\zeta = \begin{bmatrix} x(0) \\ \xi_x(1) \\ \vdots \\ \xi_x(N-1) \\ \xi_y(0) \\ \vdots \\ \xi_y(N) \end{bmatrix}.$$

We may now rewrite the equations for X_a and X_b as follows:

$$\begin{bmatrix} X_a \\ X_b \end{bmatrix} = \tilde{F}\zeta + f(u(0), \dots, u(N)).$$

Exercise 4.2 Specify \tilde{F} as well as $\tilde{f}(u)$ explicitly.

We further have,

$$\Sigma_{\zeta} := Cov(\zeta) = \begin{bmatrix} \Sigma_{x(0)} & 0 & 0 & 0 \\ & \begin{bmatrix} \Sigma_{x} & 0 \\ & \ddots & \\ 0 & \begin{bmatrix} \Sigma_{xy} & 0 \\ & \ddots & \\ 0 & \begin{bmatrix} \Sigma'_{xy} & 0 \\ & \ddots & \\ 0 & & \Sigma'_{xy} \end{bmatrix} \begin{bmatrix} \Sigma_{xy} & 0 \\ & \ddots & \\ \begin{bmatrix} \Sigma_{y} & 0 \\ & \ddots & \\ 0 & & \Sigma'_{xy} \end{bmatrix} \end{bmatrix} .$$

Here the $\Sigma_{x(0)}$ is an assumed covariance matrix for x(0). The other Σ elements are the covariances of the noise vectors: Σ_x is the covariance matrix of the disturbance. Σ_y is the covariance matrix of the measurement noise. And $\Sigma_{x,y}$ is the cross-covariance between disturbance and measurements (this is often assumed 0).

Thus,

$$Cov\left(\left[\begin{array}{c}X_a\\X_b\end{array}\right]\right) = \tilde{F}\Sigma_{\zeta}\tilde{F}' :=: \left[\begin{array}{cc}\Sigma_a & \Sigma_{ab}\\\Sigma'_{ab} & \Sigma_b\end{array}\right].$$

Observe also that,

$$\mu_a = \mathbb{E} [X_a] = [\mathbb{E} [x(0)]' \ 0' \ \dots \ 0']' + \tilde{b} (u(0), \dots, u(N-1)),$$

$$\mu_b = \mathbb{E} [X_b] = \tilde{d} (u(0), \dots, u(N)).$$

We now have all of the needed ingredients of Proposition 3.3 to calculate the LMMSE of X_a based on X_a . I.e. take,

$$h^*(x_b) = \mu_a + \Sigma_{ab} \Sigma_b^{-1}(x_b - \mu_b).$$

and then the predictor at for the state at time n is:

$$\hat{x}(n) = \left[h^*(x_b)\right]_{(nN+1,\dots,nN+n)}$$

While this is very nice, it is not efficient from a control theory perspective since getting an estimate for X_a requires computation of the order of $O((nN)^3)$. It would be much better to have some sort of recursive solution that yields $\hat{x}(N)$ at each step. This is the celebrated Kalman filtering algorithm which we present in the next section.

Exercise 4.3 Consider the scalar system:

$$\begin{aligned} x(n+1) &= 2x(n) + u(n) + \xi_x(n) \\ y(n) &= x(n) + \xi_y(n) \end{aligned}$$

Where $\xi_x(n)$ and $\xi_y(n)$ are both of unit variance and assumed uncorrelated.

Assume x(0) is such that $\mathbb{E}[x(0)] = 0$ and Var(x(0)) = 0. Assume a control input of u(n) = 1 was applied to the system over the times n = 0, 1, 2. And the measured output was, $(y(0), y(1), y(2), y(3)) = (y_0, y_1, y_2, y_3)$.

Use the derived LMMSE in this section to obtain an estimator for x(3) in terms of (y_0, y_1, y_2, y_3) .

5 The Kalman Filtering Algorithm

For simplicity in this section, we assume B = 0 and D = 0 and thus our system is

$$\begin{array}{rcl} x(n+1) &=& Ax(n) + \xi_x(n) \\ y(n) &=& Cx(n) + \xi_y(n) \end{array}$$

The more general case (with inputs) easily follows and is left as an exercise. We shall also assume for simplicity that $\Sigma_{xy} = 0$. This assumption can also be relaxed.

In general the Kalman filtering algorithm is based on (deterministic) sequence $K(0), K(1), \ldots$ that is used as follows:

$$\hat{x}(n+1) = A\hat{x}(n) + K(n)(y(n+1) - CA\hat{x}(n)).$$
(7)

In this sense it is like a Luenberger observer yet where the matrices K generally depend on time (even in the case presented here where A and C are constant). As an aid for calculating K(n) we have,

$$S(n) := Cov(x(n+1) - \hat{x}(n+1) | x(n), x(n-1), \dots, x(0)),$$

with S(n) following the following recursion:

$$S(n+1) = A\Big(S(n) - S(n)C'\big(CS(n)C' + \Sigma_y\big)^{-1}CS(n)\Big)A' + \Sigma_x.$$

Now S(n) is used to obtain K(n) as follows:

$$K(n) = S(n)C' (CS(n)C' + \Sigma_y)^{-1}.$$

Note that in many applications we may also use the steady state Kalman filter in which we take S(n) as the fixed unique positive definite S solving equation:

$$S = A \left(S - SC' \left(CSC' + \Sigma_y \right)^{-1} CS \right) A' + \Sigma_x$$

This then yields a constant K in (7).

It is obvious that the Kalman filter and (even more) the steady state Kalman filter are computationally efficient compared to the method described in the previous section.

Exercise 5.1 Consider the scalar system,

$$\begin{aligned} x(n+1) &= \frac{4}{5}x(n) + \xi_x(n), \\ y(n) &= x(n) + \xi_y(n). \end{aligned}$$

Take, $Var(\xi_x(n)) = 9/25$ and $Var(\xi_y(n)) = 1$. Find the form of the predictor $\hat{x}_n(y)$. Find the steady state predictor. We have the following:

Theorem 5.2 The sequence defined in (7) is the LMMSE estimator of x(n).

Note that the proof below is based on the the fact the noise terms are Gaussian. In this case the LMMSE is also the optimal MSE estimator. A more general proof based on the *orthogonality principle*, based on the representation of square integrable random vectors as elements of a Hilbert space is also known but is not discussed here. In that case Gaussian assumptions are not required and (7) is still the LMMSE (yet not necessarily the best MSE estimator).

Proof

Denote
$$Y(n) = (y(0), y(1), \dots, y(n))$$
 and set,
 $\hat{x}^{-}(n) := \mathbb{E} [x(n)|Y(n-1)], \qquad \hat{x}(n) := \mathbb{E} [x(n)|Y(n)].$

Observe by (4) that $\hat{x}(n)$ is the optimal MSE estimator of x(n) and thus also the LMMSE estimator since $x(\cdot)$ is Gaussian. Denote the respective conditional covariance matrices:

$$P(n) := \mathbb{E} \left[\left(x(n) - \hat{x}(n) \right) \left(x(n) - \hat{x}(n) \right)' \mid Y(n) \right], P^{-}(n) := \mathbb{E} \left[\left(x(n) - \hat{x}^{-}(n) \right) \left(x(n) - \hat{x}^{-}(n) \right)' \mid Y(n-1) \right].$$

Further for n = 0 set, $P^{-}(0) := \Sigma_{x(0)}$ and $\hat{x}^{-}(0) := \mathbb{E}[x(0)]$. Observe that in addition to $y(\cdot)$ and $x(\cdot)$, the sequences $\hat{x}^{-}(\cdot)$ and $\hat{x}(\cdot)$ are also jointly Gaussian since they are all generated by linear combinations of the "primitives" of the process, $\xi_x(\cdot), \xi_y(\cdot)$ and x(0) and also since $\hat{x}^{-}(\cdot)$ and $\hat{x}(\cdot)$ follow from the formula for the conditional expectation in (2).

The key step is to observe that when conditioning on Y(n-1), the distribution of [x(n)', y(n)']' is,

$$\mathcal{N}\Big(\begin{bmatrix}\hat{x}^{-}(n)\\C\hat{x}^{-}(n)\end{bmatrix},\begin{bmatrix}P^{-}(n)&P^{-}(n)C'\\CP^{-}(n)&CP^{-}(n)C'+\Sigma_{y}\end{bmatrix}\Big).$$
(8)

Noting that,

 $\hat{x}(n) = \mathbb{E} \left[x(n) \mid Y(n) \right] = \mathbb{E} \left[x(n) \mid y(n), Y(n-1) \right],$

we apply the mean and covariance formulas of (2) based on (8) with everything preconditioned on Y(n-1) to get:

$$\hat{x}(n) = \hat{x}^{-}(n) + P^{-}(n)C'(CP^{-}(n)C' + \Sigma_{y})^{-1}(y(n) - C\hat{x}^{-}(n)), \qquad (9)$$

$$P(n) = P^{-}(n) - P^{-}(n)C' (CP^{-}(n)C' + \Sigma_{y})^{-1}CP^{-}(n).$$
(10)

Now observe that,

$$\hat{x}^{-}(n+1) = \mathbb{E}\left[x(n+1)|Y(n)\right] = \mathbb{E}\left[Ax(n) + \xi_x(n)|Y(n)\right] = A\mathbb{E}\left[x(n)|Y(n)\right] = A\hat{x}(n),$$

and thus substitution in (9) for time n + 1 yields,

$$\hat{x}(n+1) = A\hat{x}(n) + P^{-}(n+1)C' (CP^{-}(n+1)C' + \Sigma_{y})^{-1} (y(n+1) - CA\hat{x}(n)).$$

Further,

$$P^{-}(n+1) = Cov(x(n+1) | Y(n)) = Cov(Ax(n) + \xi_x(n) | Y(n)) = AP(n)A' + \Sigma_x.$$

Substitution of (10) in the above yields

$$P^{-}(n+1) = A \Big(P^{-}(n) - P^{-}(n)C' \big(CP^{-}(n)C' + \Sigma_{y} \big)^{-1} CP^{-}(n) \Big) A' + \Sigma_{x}.$$

Now denote $S(n) := P^{-}(n+1)$ to obtain the desired equations:

$$\hat{x}(n+1) = A\hat{x}(n) + K(n)(y(n+1) - CA\hat{x}(n))
K(n) = S(n)C'(CS(n)C' + \Sigma_y)^{-1}
S(n+1) = A(S(n) - S(n)C'(CS(n)C' + \Sigma_y)^{-1}CS(n))A' + \Sigma_x.$$

Exercise 5.3 What is the Kalman filter for the case of $B \neq 0$ and $D \neq 0$. Describe any needed changes in the proof above.

6 Brief Overview of LQG

We only touch LQG briefly and informally. Consider the system,

$$\begin{aligned} x(n+1) &= Ax(n) + Bu(n) + \xi_x(n) \\ y(n) &= Cx(n) + Du(n) + \xi_y(n), \end{aligned}$$

and assume our goal is to find an optimal output feedback law: $u^*(y)$, such that the following is minimized:

$$\mathbb{E} \left[\sum_{k=0}^{N} x(n)' Q x(n) + u(n)' R u(n) \right],$$

with N either finite or infinite and $Q \ge 0$, R > 0. Assume further that (A, B) is controllable and (A, C) is observable.

This generalization of the linear quadratic regulator (LQR) problem studied in previous units, is often refereed to as the LQG problem (Linear quadratic Gaussian). Note that the LQR formulation that we studied ignored the output y and assumed state-feedback.

It turns out that solution of the LQG problem by means of dynamic programming (yet with a stochastic element) is essentially equivalent to dynamic programming solution of LQR. The basic ingredient is once again *Bellman's principle of optimality*, yet this time presented in a stochastic (Markovian) setting:

In somewhat greater generality, consider systems of the form:

$$x(n+1) = f(x(n), u(x(n)), \xi(n)), \qquad n = 0, 1, \dots, N-1,$$

where $f(\cdot)$ is some function and ξ is an i.i.d. sequence. For any prescribed $u(\cdot)$ such a system is a Markov chain (informally a stochastic process whose next step only depends on the current state and some noise component and not on the past). The basic setting of *stochastic dynamic programming* (a.k.a. Markov decision processes) is to find a $u_n^*(x), n = 0, 1, \ldots, N - 1$ such that,

$$\mathbb{E}\left[g_N(x(N)) + \sum_{k=0}^{N-1} g_k(x(k), u_k(x(k)), \xi(k))\right],$$

is minimized. Here $g_k(\cdot), k = 1, \ldots, N-1$ is the cost per stage and $g_N(\cdot)$ is the terminal cost. Note also the slight change of notation, where we put the time index as a subscript of u.

Principle of optimality (stochastic version): Let $u^* = (u_0^*(\cdot), \ldots, u_{N-1}^*(\cdot))$ be an optimal policy. Assume that in the stochastic process resulting from $u^*(\cdot)$ it is possible to reach a given state at time n. Consider now the *subproblem* whereby the process is in state x(n) at time n and wish to minimize:

$$\mathbb{E}\left[g_N(x(N)) + \sum_{k=n}^{N-1} g_k(x(k), u_k(x(k)), \xi(k))\right],$$

then the truncated policy $(u_n^*(\cdot), u_{n+1}^*(\cdot), \ldots, u_{N-1}^*(\cdot))$ is optimal for this subproblem. \Box

By application of the principle of optimality in similar spirit to as is done for the solution of discrete time LQR, we get a solution to the LQG problem that parallels that of the LQR problem, yet takes the noise into account in the following beautiful manner:

- 1. The Kalman filtering solution yields an estimator of $\hat{x}(\cdot)$.
- 2. The deterministic LQR solution (assuming known x) is applied to \hat{x} .

In view of the brevity of this section, we omit details, yet mention that this is a stochastic manifestation of the *separation principle* presented in Unit 4, where the observer and feedback control law can be designed separately and then combined. Non-linear (deterministic and stochastic) systems usually do not exhibit this clean property – and are a current active area of research.