

# Gaussian RVs and the Kalman filter

ROB 501

Necmiye Ozay

- Gaussian random variables
  - MVE another look (interpretation based on conditional probability using Schur complements)
  - (if time) Derivation of Kalman filter and why it is a minimum variance estimator

# Heart of the recursion in Kalman filter

**Definition of Terms:**

- $\hat{x}_{k|k} := \mathcal{E}\{x_k | y_0, \dots, y_k\}$  ← conditional mean at time  $k$
- $P_{k|k} := \mathcal{E}\{(x_k - \hat{x}_{k|k})(x_k - \hat{x}_{k|k})^\top | y_0, \dots, y_k\}$  ← conditional cov at time  $k$
- $\hat{x}_{k+1|k} := \mathcal{E}\{x_{k+1} | y_0, \dots, y_k\} = \mathcal{E}\{A_k x_k + G_k w_k | y_0, \dots, y_k\}$
- $P_{k+1|k} := \mathcal{E}\{(x_{k+1} - \hat{x}_{k+1|k})(x_{k+1} - \hat{x}_{k+1|k})^\top | y_0, \dots, y_k\}$

$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \sim \begin{bmatrix} x_k \\ y_k \\ z_{0, \dots, k-1} \end{bmatrix}$ ,  $X \sim x_k$ ,  $Y \sim y_k$ ,  $Z \sim \{y_0, \dots, y_{k-1}\}$

We want to compute  $X | \begin{bmatrix} Y \\ Z \end{bmatrix}$  (in Kalman filter  $\hat{x}_{k|k}$ )

→ for recursive computation, we want to express  $X | \begin{bmatrix} Y \\ Z \end{bmatrix}$  using  $X | Z$  (which is  $\hat{x}_{k|k-1}$ ) and  $Y$

$\hat{x}_{k|k-1} = A_k x_{k-1} + G_{k-1} w_{k-1}$

$\hat{x}_{k|k} = \mathcal{E}\{X | \begin{bmatrix} Y \\ Z \end{bmatrix}\}$

$X | \begin{bmatrix} Y \\ Z \end{bmatrix} \sim \frac{f_{XYZ}}{f_{YZ}}$

$\frac{f_{XYZ}}{f_Z} \sim \begin{bmatrix} X \\ Y \end{bmatrix} | Z = \begin{bmatrix} X | Z \\ Y | Z \end{bmatrix} \sim f_{(X|Z), (Y|Z)}$

$(X | Z) | (Y | Z) \sim \frac{f_{(X|Z), (Y|Z)}}{f_{(Y|Z)}}$

something we can compute recursively

$= \frac{f_{XYZ}}{f_{YZ}} = \frac{f_{XYZ}}{f_{X|Z} / f_Z} = \frac{f_{XYZ}}{f_{X|Z}} \sim X | \begin{bmatrix} Y \\ Z \end{bmatrix}$

Joint distribution of  $X, Y, Z$

Joint distribution of  $X | Z$  and  $(Y | Z)$

Joint distribution of  $X | Z$

batch version

# Rob 501 Handout: Grizzle

## Useful Facts About Gaussian Random Variables and Vectors

**Def.** A random variable  $X$  is normally distributed with mean  $\mu$  and variance  $\sigma^2 > 0$  if it has density

$$\rightarrow f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

The standard deviation is  $\sigma > 0$ . The mean and variance satisfy

$$\rightarrow \mu := \mathcal{E}\{X\} := \int_{\mathbb{R}} xf_X(x)dx := \int_{-\infty}^{\infty} xf_X(x)dx \quad \leftarrow$$

$$\rightarrow \sigma^2 := \mathcal{E}\{(X - \mu)^2\} := \int_{\mathbb{R}} (x - \mu)^2 f_X(x)dx := \int_{-\infty}^{\infty} (x - \mu)^2 f_X(x)dx \quad \leftarrow$$

**Remarks:** You should be quite familiar with the “bell curve”.  $X$  is also called a Gaussian random variable. We often say  $X$  has a *univariate normal distribution* or a *univariate Gaussian distribution* to emphasize that we are talking about a single random variable.

For the most part, we do not care too much about individual random variables. We are interested in collections of random variables and random vectors, and hence we are primarily concerned about *jointly distributed random variables*. If you take EECS 501, you can learn a tremendous amount of material about this subject. In the following, I will give a bare bones accounting of *multivariate normal random variables*.

**Def.** A finite collection random variables  $X_1, X_2, \dots, X_p$ , or equivalently, the random vector

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} \in \mathbb{R}^p$$

$$\Sigma \in \mathbb{R}^{p \times p}$$

has a (non-degenerate) *multivariate normal distribution* with mean  $\mu$  and covariance  $\Sigma > 0$  if the joint density is given by

$$f_X(x) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}.$$

joint density of  
a multivariate  
Gaussian

**Remark:** In the above,  $|\Sigma| = \det(\Sigma)$  which must be non-zero for the denominator to be well defined. This condition is what is meant by “non-degenerate”. When  $|\Sigma| = 0$ , one can still define a multivariate normal distribution, but the “moment generating function” must be used. This is a technicality that we will skip.

$$\mathcal{E}\{X\} = \mu \in \mathbb{R}^p \quad \mu_i := \int_{\mathbb{R}^p} x_i f_X(x) dx := \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} x_i f_X(x_1, \dots, x_p) dx_1 \cdots dx_p$$

$$\text{cov}(X, X) = \mathcal{E}\{(X - \mu)(X - \mu)^\top\} = \Sigma \in \mathbb{R}^{p \times p}$$

$$\mathcal{E}\{(X_i - \mu_i)(X_j - \mu_j)\} = [\Sigma]_{ij} =: \Sigma_{ij} := \int_{\mathbb{R}^p} (x_i - \mu_i)(x_j - \mu_j) f_X(x) dx$$

$$x = (x_1, x_2, \dots, x_p) = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \quad (\text{depending on context})$$

**Marginal Distributions:** Each random variable  $X_i$  has a *univariate normal distribution* with mean  $\mu_i$  and variance  $\Sigma_{ii}$ .

$$f_{X_i}(x_i) = \frac{1}{\sqrt{2\pi\Sigma_{ii}}} e^{-\frac{(x_i - \mu_i)^2}{2\Sigma_{ii}}}. \quad \checkmark$$

**Remark:** We note the unfortunate lack of coordination of notation in that the standard deviation of  $X_i$ , which we typically denote by  $\sigma_i$ , is given by

$$\sigma_i = \sqrt{\Sigma_{ii}}.$$

I guess we will not be denoting the entries of  $\Sigma$  with lower case  $\sigma$ .

**Independence:** Gaussian random variables are very special in that they are independent if, and only if, they are uncorrelated. Hence,  $X_i$  and  $X_j$  are independent if, and only if,  $\Sigma_{ij} = \Sigma_{ji} = 0$ . ✓

↳ **Linear Combinations:** Define a new random vector by  $Y = AX + b$ . Then  $Y$  is a Gaussian (normal) random vector with

$$\mathbb{E}(Y) = \mathbb{E}(AX + b) = \mathbb{E}(AX) + \underbrace{\mathbb{E}(b)}_b = A\mathbb{E}(X) + b = A\mu + b \quad \text{↳ mean of } X.$$

$$\mathbb{E}\{Y\} = A\mu + b =: \mu_Y$$

$$Y - \mu_Y = AX + b - (A\mu + b) \\ = A(X - \mu)$$

$$\mathbb{E}(A(X - \mu)(X - \mu)^T A^T) = A \underbrace{\mathbb{E}((X - \mu)(X - \mu)^T)}_{A \Sigma A^T} A^T$$

$$\text{cov}(Y, Y) = \mathbb{E}\{(Y - \mu_Y)(Y - \mu_Y)^T\} = A\Sigma A^T =: \Sigma_{YY}$$

Indeed,  $Y - \mu_Y = A(X - \mu)$ . Hence,

$$\text{cov}(Y, Y) = \mathbb{E}\{[A(X - \mu)][A(X - \mu)]^T\} = A\mathbb{E}\{(X - \mu)(X - \mu)^T\}A^T = A\Sigma A^T.$$

**Remark:** Taking  $b = 0$  and  $A$  to be a row vector with all zeros except a one in the  $i$ -th spot, that is  $A = [0, \dots, 1, \dots, 0]$ , recovers the *marginal distributions* discussed above.

$$X_i = \underbrace{[1 \ 0 \ \dots \ 0]}_A X + \underbrace{0}_b$$

**Working with Two Vectors of Gaussian Random Variables:** In addition to looking at individual random variables making up a random vector, we can group the components to form two or more blocks of vectors as long as their sizes add up to  $p$ , the number of components in  $X$ . We abuse notation and write

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \in \mathbb{R}^n$$

In books, you'll often see the blocks expressed in bold font, such as  $\mathbf{X}_1$  and  $\mathbf{X}_2$ . We will NOT do this. Conformally with this decomposition of  $X$  into two blocks, we decompose the mean and covariance as follows

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \leftarrow \quad \begin{array}{l} \mu_1 \in \mathbb{R}^n \\ \mu_2 \in \mathbb{R}^m \end{array}$$

and

$$\Sigma = \begin{bmatrix} \overset{n}{\underset{\curvearrowleft}{\Sigma_{11}}} & \overset{m}{\underset{\curvearrowleft}{\Sigma_{12}}} \\ \overset{m}{\underset{\curvearrowright}{\Sigma_{21}}} & \Sigma_{22} \end{bmatrix} \quad \begin{array}{c} \overset{n}{\underset{\curvearrowleft}{\curvearrowright}} \\ \overset{m}{\underset{\curvearrowright}{\curvearrowleft}} \end{array}$$

**Remark:** From our results on the Schur complement, we know that  $\Sigma > 0$  if, and only if,  $\Sigma_{22} > 0$  and  $\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} > 0$ .

To be super clear on the dimensions, we suppose  $n + m = p$  and note that

$$\mu_1 = \mathcal{E}\{X_1\} \in \mathbb{R}^n \quad \text{and} \quad \mu_2 = \mathcal{E}\{X_2\} \in \mathbb{R}^m$$

$$\text{cov}(X_1, X_1) = \Sigma_{11} \in \mathbb{R}^{n \times n} \quad \text{and} \quad \text{cov}(X_2, X_2) = \Sigma_{22} \in \mathbb{R}^{m \times m}$$

$$\text{cov}(X_1, X_2) = \Sigma_{12} \in \mathbb{R}^{n \times m} \quad \text{and} \quad \text{cov}(X_2, X_1) = \Sigma_{21} \in \mathbb{R}^{m \times n}$$

Furthermore, because  $\Sigma = \Sigma^\top$ , we have that

$$\Sigma_{11}^\top = \Sigma_{11}, \quad \Sigma_{22}^\top = \Sigma_{22}, \quad \text{and} \quad \Sigma_{12}^\top = \Sigma_{21}.$$

**Remark:** Each vector  $X_i$  has a multivariate normal distribution with mean  $\mu_i$  and covariance  $\Sigma_{ii}$ . This is also called the **marginal distribution** of  $X_i$ . If we know the mean and covariance for the composite vector  $X$ , it is very easy to read off the marginal distributions of its vector components.

Gaussian random variable

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim N \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

**Fact 1: Conditional Distributions of Gaussian Random Vectors:** Let  $X_1$  and  $X_2$  be as above, namely they are components of a larger vector  $X$  that has a multivariate normal distribution. Then the conditional distribution of  $X_1$  given  $X_2 = x_2$  has a multivariate normal distribution with

Mean :  $\mu_{1|2} := \mu_1 + \underbrace{\Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2)}$   
 and Covariance:  $\Sigma_{1|2} := \Sigma_{11} - \underbrace{\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}}$ .

In passing, we note that the mean depends on the value of  $x_2$  while the covariance does not.

To be extra clear on the meanings here,

- $\mu_{1|2} = \mathcal{E}\{X_1 \mid X_2 = x_2\}$
- $\Sigma_{1|2} = \mathcal{E}\{(X_1 - \mu_{1|2})(X_1 - \mu_{1|2})^\top \mid X_2 = x_2\}$
- $X_1$  given  $X_2 = x_2$  is a random vector. It has a multivariate normal distribution with the above mean vector and covariance matrix. Specifically, its density is

$$f_{X_1|X_2=x_2}(x_1) = \frac{f_{X_1, X_2}(x_1, x_2)}{f_{X_2}(x_2)} = \frac{1}{\sqrt{(2\pi)^n |\Sigma_{1|2}|}} e^{-\frac{1}{2}(x_1 - \mu_{1|2})^\top \Sigma_{1|2}^{-1}(x_1 - \mu_{1|2})},$$

where it is emphasized that  $\mu_{1|2}$  depends explicitly on  $x_2$ .

**Remark:** A proof of this can be found at the link below. The algebra is rather painful. If you are very ambitious, you can work out the special case where  $X_1$  and  $X_2$  are scalars. This will not be on any exam.

<http://fourier.eng.hmc.edu/e161/lectures/gaussianprocess/node7.html>  
 See also

<http://www.stats.ox.ac.uk/~steffen/teaching/bs2HT9/gauss.pdf>.

**Fact 2 on Conditional Independence:** Suppose we have 3 vectors  $X_1$ ,  $X_2$  and  $X_3$  that are jointly normally distributed:

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix}, \Sigma\right)$$

and that  $X_2$  is independent of  $X_1$  and  $X_3$ . We then have no special structure on the means,

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix}$$

but the covariance matrix has the form

$$\Sigma = \begin{bmatrix} \Sigma_{11} & 0 & \Sigma_{13} \\ 0 & \Sigma_{22} & 0 \\ \Sigma_{13}^\top & 0 & \Sigma_{33} \end{bmatrix}$$

where  $\Sigma_{12} = \Sigma_{21}^\top = \text{cov}(X_1, X_2) = 0$  due to the independence of  $X_1$  and  $X_2$ . Similarly for  $\Sigma_{23} = \Sigma_{32}^\top = 0$ . Because  $\Sigma$  is symmetric,  $\Sigma_{31} = \Sigma_{13}^\top$ .

**Then  $X_1$  and  $X_2$  are conditionally independent given  $X_3$ .** Written a different way, the two normal random variables,  $X_{1|X_3}$  ( $X_1$  conditioned on knowing  $X_3$ ) and  $X_{2|X_3}$  ( $X_2$  conditioned on knowing  $X_3$ ) are independent.

$$\text{Cov}((X_1|X_3), (X_2|X_3)) = 0$$

high-level intuition:  $X_2|X_3 \leftrightarrow X_2$  (bec.  $X_2$  and  $X_3$  are independent)  
 and  $(X_1|X_3)$  are independent.

To see why this is true, we partition  $\Sigma$  as

$$\Sigma = \begin{bmatrix} \Sigma_{11} & 0 & \Sigma_{13} \\ 0 & \Sigma_{22} & 0 \\ \Sigma_{13} & 0 & \Sigma_{33} \end{bmatrix}.$$

We compute the covariance of  $X_1$  and  $X_2$  conditioned on  $X_3$ , that is

$$\left[ \begin{array}{c} X_1 \\ X_2 \end{array} \right] \Big| X_3, = \left[ \begin{array}{c} X_1 | X_3 \\ X_2 | X_3 \end{array} \right]$$

using the Schur complement from **Fact 1**

$$\begin{aligned} \text{cov}\left(\left[ \begin{array}{c} X_1 | X_3 \\ X_2 | X_3 \end{array} \right], \left[ \begin{array}{c} X_1 | X_3 \\ X_2 | X_3 \end{array} \right]\right) &= \left[ \begin{array}{cc} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} \end{array} \right] - \left[ \begin{array}{c} \Sigma_{13} \\ 0 \end{array} \right] \Sigma_{33}^{-1} \left[ \begin{array}{cc} \Sigma_{13}^\top & 0 \end{array} \right] \\ &= \left[ \begin{array}{cc} \Sigma_{11} - \Sigma_{13} \Sigma_{33}^{-1} \Sigma_{13}^\top & 0 \\ 0 & \Sigma_{22} \end{array} \right] \end{aligned}$$

$$\Sigma_{12} := \underbrace{\Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}}_{\sim}$$

Because the off-diagonal blocks are zero, the two random variables  $X_1|X_3$  and  $X_2|X_3$  are uncorrelated, and because they are normal, we conclude they are independent.

Once again, what we have seen is that if  $X_1$  and  $X_2$  are independent, and we also have  $X_2$  is independent of  $X_3$ , then  $X_1$  and  $X_2$  remain independent when we condition them on  $X_3$ .

**Fact 3: Covariance of a Sum of Independent Normal Random Variables** Let  $X_1$  and  $X_2$  be independent normal random vectors, with means  $\mu_1$  and  $\mu_2$ , and covariances,  $\Sigma_{11}$  and  $\Sigma_{22}$ . Define  $Y$  as a “linear combination”  $X_1$  and  $X_2$  via

$$\rightarrow Y = \underbrace{AX_1 + BX_2}_{\text{for appropriately sized matrices } A \text{ and } B. \text{ Then}}$$

$$\begin{aligned} E(AX_1 + BX_2) &= E(AX_1) + E(BX_2) \\ &= A\mu_1 + B\mu_2 \end{aligned}$$

and

$$\text{cov}(Y, Y) = A\Sigma_{11}A^\top + B\Sigma_{22}B^\top$$

$$AX_1 + BX_2 - A\mu_1 - B\mu_2 = A(X_1 - \mu_1) - B(X_2 - \mu_2)$$

To see why this is true, we first note that

$$\begin{aligned} (Y - \mu_Y)(Y - \mu_Y)^\top &= A(X_1 - \mu_1)(X_1 - \mu_1)^\top A^\top + B(X_2 - \mu_2)(X_2 - \mu_2)^\top B^\top \\ &\quad + 2A(X_1 - \mu_1)(X_2 - \mu_2)^\top B^\top \end{aligned}$$

And then note that when expectations are taken on each side, the independence of  $X_1$  and  $X_2$  gives

$$E\{(X_1 - \mu_1)(X_2 - \mu_2)^\top\} = 0.$$

Therefore,

$$\begin{aligned} \text{cov}(Y, Y) &= E\{(Y - \mu_Y)(Y - \mu_Y)^\top\} \\ &= A E\{(X_1 - \mu_1)(X_1 - \mu_1)^\top\} A^\top + B E\{(X_2 - \mu_2)(X_2 - \mu_2)^\top\} B^\top \\ &= A\Sigma_{11}A^\top + B\Sigma_{22}B^\top. \end{aligned}$$

**Remark:** The next few pages discuss the “information” or “precision” matrix. You will likely encounter it in other courses, such as Mobile Robotics, or in papers. We first saw it when we compared BLUE to Weighted Least Squares (they are the same when the weight chosen as the Information Matrix of the

noise term). You will not need to know anything about the information matrix in the context of the Kalman Filter: when seeing the filter for the first time, you do not need to do every possible variation.

**Information or Precision Matrix:** The Kalman filter can be written in many forms. One alternative form propagates the inverse of the covariance matrix instead of the covariance matrix. The inverse of the covariance matrix has two names: *information matrix* and *precision matrix*. We will use the first one:

$$\text{Information matrix: } \Lambda := \Sigma^{-1}$$

We decompose it just as we did with the covariance matrix.

$$\Lambda =: \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix}$$

The formula for inversion of block matrices gives

$$\begin{aligned}\Lambda_{11} &= (\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1} \\ \Lambda_{12} &= -(\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1}\Sigma_{12}\Sigma_{22}^{-1} \\ \Lambda_{21} &= \Lambda_{12}^\top \\ \Lambda_{22} &= \Sigma_{22}^{-1} + \Sigma_{22}^{-1}\Sigma_{21}(\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1}\Sigma_{12}\Sigma_{22}^{-1}\end{aligned}$$

(See [http://en.wikipedia.org/wiki/Matrix\\_inversion\\_lemma#Blockwise\\_inversion](http://en.wikipedia.org/wiki/Matrix_inversion_lemma#Blockwise_inversion))

We also scale the mean by defining

$$\eta := \Lambda\mu$$

that is,

$$\begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} := \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

**Remark:** For a multivariate normal distribution, it is equivalent to know  $\eta$  and  $\Lambda$  or  $\mu$  and  $\Sigma$ . We go back and forth between the two by matrix inversion and multiplication. One sometimes says that these are dual parameterizations for the normal distribution. We only mention the alternative parameterization with the information matrix because sometimes it is easier to use than the more standard mean and covariance representation.

**Conditional Distributions Using the Information Matrix:** The information matrix of the random variable  $X_1$  given that  $X_2 = x_2$  is

$$\Lambda_{1|2} = \Lambda_{11}$$

and

$$\eta_{1|2} = \eta_1 - \Lambda_{12}x_2$$

In other words, if you have the information matrix handy, computing the conditional distribution is easier with it than with the covariance matrix. We note that if you want to go back to the standard representation, then

$$\Sigma_{1|2} = \Lambda_{1|2}^{-1}$$

and

$$\mu_{1|2} = \Lambda_{1|2}^{-1} \eta_{1|2}$$

**Marginal Distributions Using the Information Matrix:** Getting the marginal distributions from the information form of the distribution is more complicated. If you are interested, you can easily find it on the web or in most graduate level probability texts.

Revisiting MVE: for the special case

of normal R.V.s (also keep in mind  
that MVE is a special case of  
Kalman filter)

$$y = Cx + \varepsilon, \quad x \in \mathbb{R}^n \quad x \sim N(\mu_x, \Sigma_{xx}) \\ \varepsilon \in \mathbb{R}^m \quad \varepsilon \sim N(\mu_\varepsilon, \Sigma_{\varepsilon\varepsilon})$$

$$\mathbb{E}\{(x - \mu_x)(\varepsilon - \mu_\varepsilon)^\top\} = \Sigma_{x\varepsilon} = \Sigma_{\varepsilon x}^\top = 0_{n \times m}$$

By Fact 3:

$$\mu_y = \mathbb{E}(Cx + \varepsilon) = C\mu_x + \mu_\varepsilon$$

$$\Sigma_{yy} = \mathbb{E}\{(y - \mu_y)(y - \mu_y)^\top\}$$

$$= C\Sigma_{xx}C^\top + \Sigma_{\varepsilon\varepsilon}$$

$$\underline{V} = \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\mathbb{E}(V) = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}$$

$$\text{Cov}(V) = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx}^\top & \Sigma_{yy} \end{bmatrix}$$

$$\boxed{\begin{array}{l} \text{MVE } (\mu_x = 0, \mu_\varepsilon = 0) \\ \Sigma_{\varepsilon\varepsilon} = Q, \Sigma_{xx} = P \\ \hat{x} = P C^\top (C P C^\top + Q)^{-1} y \end{array}}$$

$$Y = A X_1 + B X_2 \\ \text{cov}(Y, Y) = A \Sigma_{11} A^\top + B \Sigma_{22} B^\top$$

$$Y = C X + I \varepsilon$$

$$C x + \varepsilon - \mu_y \\ = C x + \varepsilon - C \mu_x - \mu_\varepsilon$$

$$\begin{aligned}
\Sigma_{xy} &= E \left\{ (x - \mu_x)(y - \mu_y)^T \right\} \\
&= E \left\{ (x - \mu_x)(Cx + \varepsilon - \mu_y)^T \right\} \\
&= E \left\{ (x - \mu_x)(Cx + \varepsilon - C\mu_x - \mu_\varepsilon)^T \right\} \\
&= E \left\{ (x - \mu_x)(C(x - \mu_x) + \varepsilon - \mu_\varepsilon)^T \right\} \\
&= E \left\{ (x - \mu_x)(x - \mu_x)^T C^T \right\} + E \left\{ (x - \mu_x)(\varepsilon - \mu_\varepsilon)^T \right\} \\
&= \Sigma_{xx} C^T + 0
\end{aligned}$$

$$\text{Cov}(V) = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xx} C^T \\ C \Sigma_{xx} & C \Sigma_{xx} C^T + \Sigma_{\varepsilon\varepsilon} \end{bmatrix}$$

Let's look at  $X | Y$ , Mean :  $\mu_{1|2} := \mu_1 + \underline{\Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2)}$

$$\text{Covariance: } \Sigma_{1|2} := \underline{\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}}$$

From Fact 1:

$$\begin{aligned}
\mu_{x|y} &= \underline{\mu_x + \Sigma_{xx} C^T (C \Sigma_{xx} C^T + \Sigma_{\varepsilon\varepsilon})^{-1} (y - \mu_y)} \\
\Sigma_{x|y} &= \underline{\Sigma_{xx} - \Sigma_{xx} C^T (C \Sigma_{xx} C^T + \Sigma_{\varepsilon\varepsilon})^{-1} C \Sigma_{xx}}
\end{aligned}$$

$$\hat{x} = E(x|y) \quad \text{with} \\ \mu_x = 0 \\ \mu_{\varepsilon} = 0$$

$$= \mu_{x|y} \\ = P C^T (C P C^T + Q)^{-1} y$$

NVE ( $\mu_x = 0, \mu_{\varepsilon} = 0$   
 $\Sigma_{\varepsilon\varepsilon} = Q, \underline{\Sigma_{xx} = P}$ )

$$\hat{x} = P C^T (C P C^T + Q)^{-1} y$$

Measurement

$BX$

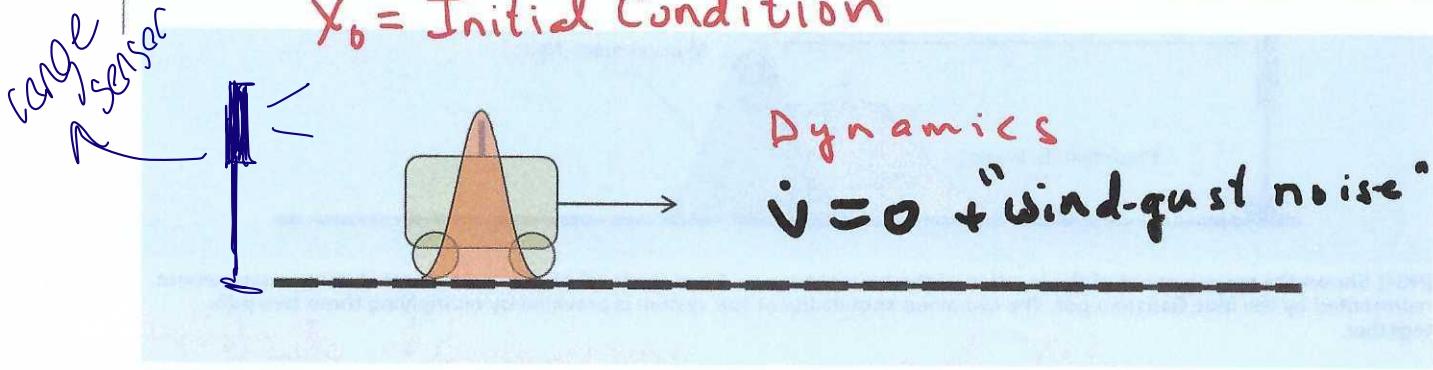
MVE

$$x_{t+1} = I x_t + Q w_t \\ y_t = C x_t + I \varepsilon_t$$

# Kalman Filter Motivation

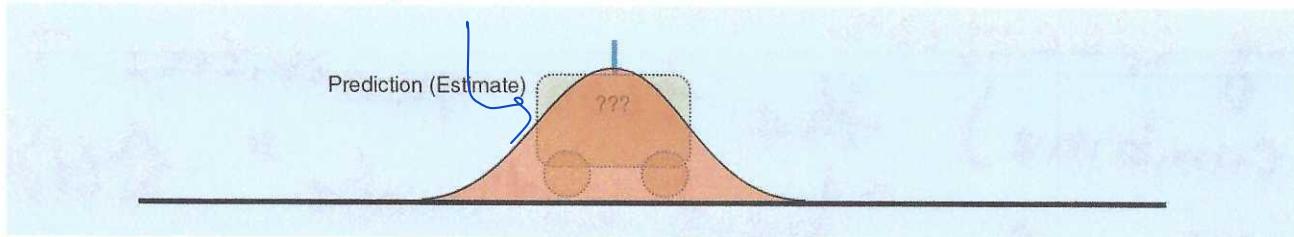
lecture **NOTES** continued

$x_0 = \text{Initial Condition}$

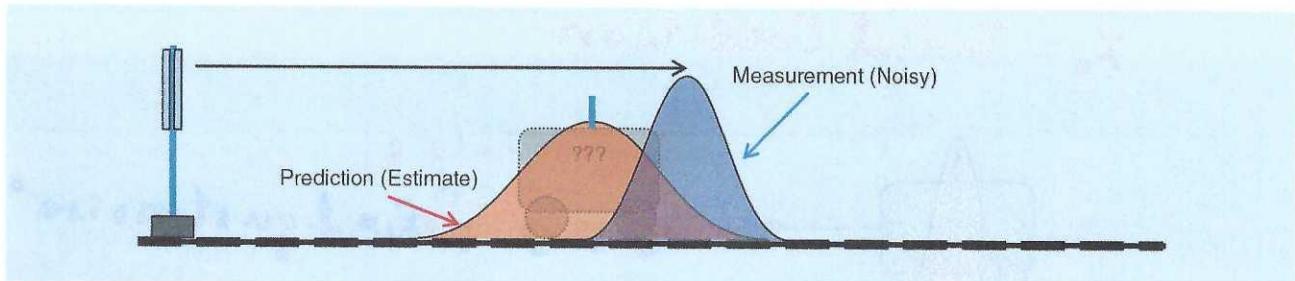


[FIG2] The initial knowledge of the system at time  $t = 0$ . The red Gaussian distribution represents the initial uncertainty.

Uncertainty typically grows as time increases

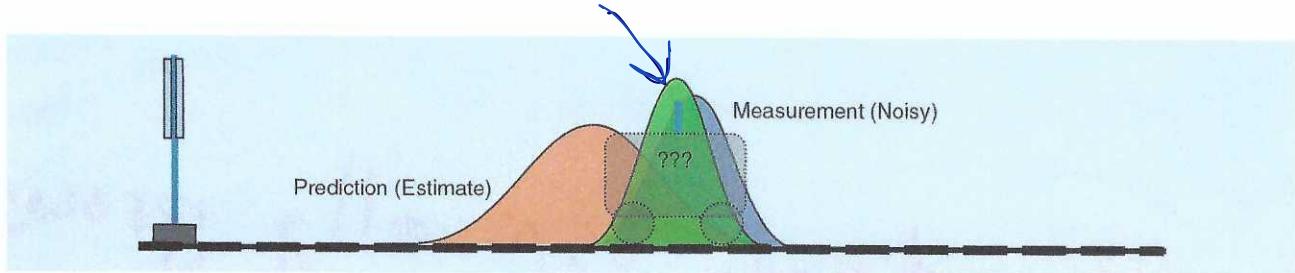


[FIG3] Here, the prediction of the location of the train at time  $t = 1$  and the level of uncertainty in that prediction is shown. The confidence in the knowledge of the position of the train has decreased, as we are not certain if the train has undergone any accelerations or decelerations in the intervening period from  $t = 0$  to  $t = 1$ .



[FIG4] Shows the measurement of the location of the train at time  $t = 1$  and the level of uncertainty in that noisy measurement, represented by the blue Gaussian pdf. The combined knowledge of this system is provided by multiplying these two pdfs together.

$x(t)$  has a "probability distribution"  
 $y(t)$  is not perfect = "probability distribution"



[FIG5] Shows the new pdf (green) generated by multiplying the pdfs associated with the prediction and measurement of the train's location at time  $t = 1$ . This new pdf provides the best estimate of the location of the train, by fusing the data from the prediction and the measurement.

Big Question: How to "fuse" (= combine) the two quantities to have a "best estimate"  $\hat{x}(t)???$

$$\text{RLS: } A_k = I, \quad G_k = 0$$

$V_k$  = deterministic =  
 = no stochastic or random model.

# Rob 501 Handout: Grizzle

## Discrete-Time Kalman Filter and its Derivation

**Model** Linear time-varying discrete-time system with “white<sup>1</sup>” Gaussian noise

$$\begin{aligned}x_{k+1} &= A_k x_k + G_k w_k, \quad x_0 \text{ initial condition} \\y_k &= C_k x_k + v_k\end{aligned}$$

$x \in \mathbb{R}^n$ ,  $w \in \mathbb{R}^p$ ,  $y \in \mathbb{R}^m$ ,  $v \in \mathbb{R}^m$ . Moreover, the random vectors  $x_0$ , and, for  $k \geq 0$ ,  $w_k$ ,  $v_k$  are all independent<sup>2</sup> Gaussian (normal) random vectors.

**Notation:**  $\delta_{kl} = 1 \Leftrightarrow k = l$  (and  $\delta_{kl} = 0$ ,  $k \neq l$ )

### Precise Assumptions on the Random Variables

- For all  $k \geq 0$ ,  $l \geq 0$ ,  $x_0$ ,  $w_k$ ,  $v_l$  are jointly Gaussian.
- $w_k$  is a 0-mean white noise process:  $\mathcal{E}\{w_k\} = 0$ , and  $\text{cov}(w_k, w_l) = R_k \delta_{kl}$
- $v_k$  is a 0-mean white noise process:  $\mathcal{E}\{v_k\} = 0$ , and  $\text{cov}(v_k, v_l) = Q_k \delta_{kl}$
- Uncorrelated noise processes:  $\text{cov}(w_k, v_l) = 0$
- The initial condition  $x_0$  is uncorrelated with all other noise sequences.
- We denote the mean and covariance of  $x_0$  by

$$\bar{x}_0 = \mathcal{E}\{x_0\} \quad \text{and} \quad P_0 = \text{cov}(x_0) = \text{cov}(x_0, x_0) = \mathcal{E}\{(x_0 - \bar{x}_0)(x_0 - \bar{x}_0)^\top\}$$

---

<sup>1</sup>Recall that in white light, all frequencies are present. When only certain frequency components are present, you get “colored” light, such as blue light or red light. The term “white” noise means that if you compute the power spectral density of the noise random process, it is a constant, meaning that all frequency components are equally represented, just as in white light.

<sup>2</sup>Recall that for normal random variables, uncorrelated and independent are the same thing. This is one of several special properties of Gaussian random variables.

**Short-hand notation for the noise modeling assumptions:**

$$\text{cov} \left( \begin{bmatrix} w_k \\ v_k \\ x_0 \end{bmatrix}, \begin{bmatrix} w_l \\ v_l \\ x_0 \end{bmatrix} \right) = \begin{bmatrix} R_k \delta_{kl} & 0 & 0 \\ 0 & Q_k \delta_{kl} & 0 \\ 0 & 0 & P_0 \end{bmatrix}, \quad \delta_{kl} = \begin{cases} 1, & k = l \\ 0, & k \neq l \end{cases}$$

**Lemma (Properties of  $x_k$  and  $y_k$  Coming from the Model)**

- For all  $k \geq 1$ ,  $x_k$  is a linear combination of  $x_0$  and  $w_0, \dots, w_{k-1}$ . In particular,  $x_k$  is uncorrelated with  $w_k$ .
- For all  $k \geq 1$ ,  $y_k$  is a linear combination of  $x_0, w_0, \dots, w_{k-1}$ , and  $v_0, \dots, v_k$ . In particular,  $y_k$  is uncorrelated with  $w_k$ .
- For all  $k \geq 0$ ,  $v_k$  is uncorrelated with  $x_k$ .

The proof is by induction using the recursive nature of the discrete-time model. We skip it. The reader can easily fill it in.

**Remark:** On the next page, we give (one form of) the discrete-time Kalman Filter. After that, we do the main elements of its derivation. There are many variations of the basic filter, all equivalent to the one we give, but some preferable over others for numerical reasons. Page 8 provides a version of the filter with the measurement update and prediction steps combined.

## The Kalman Filter

**Definition of Terms:**

$$\begin{aligned}\hat{x}_{k|k} &:= \mathcal{E}\{x_k | y_0, \dots, y_k\} \\ P_{k|k} &:= \mathcal{E}\{(x_k - \hat{x}_{k|k})(x_k - \hat{x}_{k|k})^\top | y_0, \dots, y_k\}\end{aligned}$$

$$\begin{aligned}\hat{x}_{k+1|k} &:= \mathcal{E}\{x_{k+1} | y_0, \dots, y_k\} \\ P_{k+1|k} &:= \mathcal{E}\{(x_{k+1} - \hat{x}_{k+1|k})(x_{k+1} - \hat{x}_{k+1|k})^\top | y_0, \dots, y_k\}\end{aligned}$$

**Initial Conditions:**

$$\hat{x}_{0|-1} := \bar{x}_0 = \mathcal{E}\{x_0\}, \quad \text{and} \quad P_{0|-1} := P_0 = \text{cov}(x_0)$$

**For**  $k \geq 0$

**Measurement Update Step:**

$$\begin{aligned}K_k &= P_{k|k-1} C_k^\top (C_k P_{k|k-1} C_k^\top + Q_k)^{-1} \\ &\quad (\text{Kalman Gain}) \\ \hat{x}_{k|k} &= \hat{x}_{k|k-1} + K_k (y_k - C_k \hat{x}_{k|k-1}) \\ P_{k|k} &= P_{k|k-1} - K_k C_k P_{k|k-1}\end{aligned}$$

**Time Update or Prediction Step:**

$$\begin{aligned}\hat{x}_{k+1|k} &= A_k \hat{x}_{k|k} \\ P_{k+1|k} &= A_k P_{k|k} A_k^\top + G_k R_k G_k^\top\end{aligned}$$

**End of For Loop** (Just stated this way to emphasize the recursive nature of the filter)

## Preliminaries

**Measurements:** We collect all of the measurements at time  $k$  as

$$Y_k = (y_k, y_{k-1}, \dots, y_0).$$

Strictly speaking, we should be stacking them up into a column vector as we have done for all of our estimation problems, but notationally, it is more convenient to write them in a row. Also, it is more convenient to put the most recent measurement at the head of the list. We note that

$$Y_k = (y_k, Y_{k-1}).$$

Hence,

$$\begin{aligned} \hat{x}_{k|k} &:= \mathcal{E}\{x_k|Y_k\} \\ P_{k|k} &:= \mathcal{E}\{(x_k - \hat{x}_{k|k})(x_k - \hat{x}_{k|k})^\top|Y_k\} \\ &\text{mean and covariance of the conditional normal random vector } x_k|Y_k \end{aligned}$$

$$\begin{aligned} \hat{x}_{k+1|k} &:= \mathcal{E}\{x_{k+1}|Y_k\} \\ P_{k+1|k} &:= \mathcal{E}\{(x_{k+1} - \hat{x}_{k+1|k})(x_{k+1} - \hat{x}_{k+1|k})^\top|Y_k\} \\ &\text{mean and covariance of the conditional normal random vector } x_{k+1}|Y_k \end{aligned}$$

### Important Remarks:

- The conditional random vector  $x_k|Y_k$  is distributed  $N(\hat{x}_{k|k}, P_{k|k})$ .
- The conditional random vector  $x_{k+1}|Y_k$  is distributed  $N(\hat{x}_{k+1|k}, P_{k+1|k})$ .

## Filter Derivation Using Induction and Properties of Conditional Distributions of Gaussian Random Vectors

**Base step:** The initial conditions of the filter at time  $k = 0$ , namely

$$\hat{x}_{0|-1} := \bar{x}_0, \quad \text{and} \quad P_{0|-1} := P_0$$

**Induction step:** At time  $k \geq 0$ , we suppose that  $(\hat{x}_{k|k-1}, P_{k|k-1})$  are known, and we derive  $(\hat{x}_{k|k}, P_{k|k})$  and  $(\hat{x}_{k+1|k}, P_{k+1|k})$ .

**Key idea of the development:** We need to compute the distribution (or density) of the conditional random vector

$$x_k | Y_k = x_k | (y_k, Y_{k-1})$$

From HW 7, we learned **Fact 4**  $X|(Y, Z) = X|Z \mid Y|Z$ . From this we obtain

$$x_k | Y_k = x_k | (y_k, Y_{k-1}) = x_k | Y_{k-1} \mid y_k | Y_{k-1} \quad (*)$$

where we have identified

$$x_k \leftrightarrow X, \quad y_k \leftrightarrow Y, \quad \text{and} \quad Y_{k-1} \leftrightarrow Z.$$

Hence, if we can compute the distribution (or density) of

$$\begin{bmatrix} x_k \\ y_k \end{bmatrix} \mid Y_{k-1}$$

then we can apply **Fact 1** to obtain (\*).

The following calculations are aimed at doing just this.

**Measurement Update:** We seek to derive the filter equations on page 3. To begin, we have that  $y_k = C_k x_k + v_k$ . It follows by linearity that the conditional random variable  $y_k|Y_{k-1}$  is equal to

$$y_k|Y_{k-1} = C_k x_k|Y_{k-1} + v_k|Y_{k-1}.$$

From our assumptions on the noise,  $v_k$  is independent of both  $x_k$  and  $Y_{k-1}$ , and hence by **Fact 2**,  $v_k|Y_{k-1}$  is independent of the conditional random variable  $x_k|Y_{k-1}$ . Moreover, because  $v_k$  is independent of  $Y_{k-1}$ ,  $v_k|Y_{k-1} = v_k$ . Putting this together, we have that

$$y_k|Y_{k-1} = C_k x_k|Y_{k-1} + v_k,$$

and  $x_k|Y_{k-1}$  and  $v_k$  are independent. Hence

$$\begin{aligned}\hat{y}_{k|k-1} &:= \mathcal{E}\{y_k|Y_{k-1}\} \\ &= \mathcal{E}\{C_k x_k|Y_{k-1}\} + \mathcal{E}\{v_k\} \\ &= C_k \mathcal{E}\{x_k|Y_{k-1}\} + \mathcal{E}\{v_k\} \\ &= C_k \hat{x}_{k|k-1} + 0 \\ &= C_k \hat{x}_{k|k-1}.\end{aligned}$$

Moreover, the independence of  $x_k|Y_{k-1}$  and  $v_k$  with **Fact 3** yields

$$\text{cov}(y_k|Y_{k-1}, y_k|Y_{k-1}) = C_k P_{k|k-1} C_k^\top + Q_k.$$

We use independence again to obtain

$$\text{cov}(x_k|Y_{k-1}, y_k|Y_{k-1}) = \text{cov}(x_k|Y_{k-1}, C_k x_k|Y_{k-1}) = P_{k|k-1} C_k^\top.$$

With this information, we conclude that the vector

$$\begin{bmatrix} x_k \\ y_k \end{bmatrix} | Y_{k-1}$$

is jointly normally distributed, with mean and covariance

$$\begin{bmatrix} \hat{x}_{k|k-1} \\ C_k \hat{x}_{k|k-1} \end{bmatrix}, \begin{bmatrix} P_{k|k-1} & P_{k|k-1} C_k^\top \\ C_k P_{k|k-1} & C_k P_{k|k-1} C_k^\top + Q_k \end{bmatrix} \quad (*)$$

As discussed on the previous page, to compute the distribution of  $(x_k | Y_k)$ , we have from **Fact 4**

$$x_k | Y_k = x_k \Big| (y_k, Y_{k-1}) = x_k | Y_{k-1} \Big| y_k | Y_{k-1},$$

and thus applying **Fact 1** to  $(*)$  we compute the mean and covariance of  $x_k | Y_k = x_k | Y_{k-1} \Big| y_k | Y_{k-1}$  to be

$$\hat{x}_{k|k} = \hat{x}_{k|k-1} + P_{k|k-1} C_k^\top [C_k P_{k|k-1} C_k^\top + Q_k]^{-1} (y_k - C_k \hat{x}_{k|k-1})$$

$$P_{k|k} = P_{k|k-1} - P_{k|k-1} C_k^\top [C_k P_{k|k-1} C_k^\top + Q_k]^{-1} C_k P_{k|k-1}$$

**Remark:** We note that  $P_{k|k}$  is the Schur complement  $C_k P_{k|k-1} C_k^\top + Q_k$  in the covariance of

$$\begin{bmatrix} x_k \\ y_k \end{bmatrix} | Y_{k-1}$$

**Prediction or Time Update:** We seek to derive the filter equations on page 3. This time we use the state-variable model instead of the output model, namely

$$x_{k+1} = A_k x_k + G_k w_k,$$

and we are interested in the random vector

$$x_{k+1}|Y_k = A_k x_k|Y_k + G_k w_k|Y_k.$$

Because  $x_k$  and  $Y_k$  are both independent of  $w_k$ , by **Fact 2**,  $x_k|Y_k$  and  $w_k|Y_k$  are also independent. It follows that

$$\begin{aligned}\widehat{x}_{k+1|k} &= \mathcal{E}\{x_{k+1}|Y_k\} \\ &= \mathcal{E}\{A_k x_k + G_k w_k|Y_k\} \\ &= A_k \mathcal{E}\{x_k|Y_k\} + G_k \mathcal{E}\{w_k|Y_k\} \\ &= A_k \widehat{x}_{k|k} + G_k \mathcal{E}\{w_k\} \\ &= A_k \widehat{x}_{k|k},\end{aligned}$$

where we have used  $w_k|Y_k = w_k$ , and  $\mathcal{E}\{w_k\} = 0$ .

Next, we use **Fact 3** and the conditional independence of the random vectors  $x_k|Y_k$  and  $w_k|Y_k$  to evaluate the covariance of  $x_{k+1}|Y_k$  as

$$P_{k+1|k} = A_k P_{k|k} A_k^\top + G_k R_k G_k^\top.$$

**That's the Proof Folks!**

## Combined Update Version of the Filter (often easier to implement)

Here, we will assume the model also has a *deterministic* input  $u_k$ , and thus

$$\begin{aligned}x_{k+1} &= A_k x_k + B_k u_k + G_k w_k \\y_k &= C_k x_k + v_k,\end{aligned}$$

with the assumptions on the random vectors  $x_0$ ,  $w_k$  and  $v_k$  the same as before.

**Combined Filter:** The measurement-update step and time-update step of the Kalman Filter can be combined into a single step. The algorithm becomes:

**Initial Conditions:**

$$\hat{x}_{0|-1} := \bar{x}_0 = \mathcal{E}\{x_0\}, \quad \text{and} \quad P_{0|-1} := P_0 = \text{cov}(x_0)$$

**For**  $k \geq 0$

$$K_k = (P_{k|k-1} C_k^\top) [C_k P_{k|k-1} C_k^\top + Q_k]^{-1}$$

$$\hat{x}_{k+1|k} = A_k \hat{x}_{k|k-1} + B_k u_k + A_k K_k (y_k - C_k \hat{x}_{k|k-1})$$

$$P_{k+1|k} = A_k [P_{k|k-1} - K_k C_k P_{k|k-1}] A_k^\top + G_k R_k G_k^\top$$

**End of For Loop**

**Remark:** You do not have to start at  $k = 0$ . In MATLAB, it is often easier to begin with  $k = 1$ . In that case, the initial conditions are

$$\hat{x}_{1|0} := \bar{x}_0 = \mathcal{E}\{x_0\}, \quad \text{and} \quad P_{1|0} := P_0 = \text{cov}(x_0)$$

**Remark:**  $K_k C_k P_{k|k-1} = (P_{k|k-1} C_k^\top) [C_k P_{k|k-1} C_k^\top + Q_k]^{-1} C_k P_{k|k-1}$  is symmetric positive semi-definite and represents the value of the measurement in reducing the covariance of the state estimate, just as in the MVE.