

# Minimum Mean-Square Error Estimation

Now we will take our first look at estimating variables that are themselves random, subject to a known probability law.

We start our discussion with a very basic problem. Suppose  $Y$  is a scalar random variable with a known pdf  $f_Y(y)$ . Here is a fun game: you guess what  $Y$  is going to be, then I draw a realization of  $Y$  corresponding to its probability law, then we see how close you were with your guess.

What is your best guess?

Well, that of course depends on what exactly we mean by “best”, i.e. what price I pay for being a certain amount off. But if we penalize the *mean-squared error*, we know exactly how to minimize it.

Let  $g$  be your guess. The error in your guess is of course random (since the realization of  $Y$  is random), and so is the squared-error  $(Y - g)^2$ . We want to choose  $g$  so that the mean of the squared error is as small as possible:

$$\underset{g}{\text{minimize}} \quad \mathbb{E}[(Y - g)^2].$$

Expanding the squared error makes it clear how to do this:

$$\mathbb{E}[(Y - g)^2] = \mathbb{E}[Y^2] - 2g \mathbb{E}[Y] + g^2.$$

No matter what the first moment  $\mathbb{E}[Y]$  and second moment  $\mathbb{E}[Y^2]$  are (as long as they are finite), the expression above is a convex quadratic function in  $g$ , and hence is minimized when its first derivative (w.r.t.  $g$ ) is zero, i.e. when

$$-2 \mathbb{E}[Y] + 2g = 0 \quad \Rightarrow \quad \hat{g} = \mathbb{E}[Y].$$

The squared error for this choice  $\hat{g}$  is of course exactly the variance of  $Y$ .

The story gets more interesting (and relevant) when we have multiple random variables, some of which we observe, some of which we do not. Suppose that two random variables  $(Y, Z)$  have joint pdf  $f_{Y,Z}(y, z)$ . Suppose that a realization of  $(Y, Z)$  is drawn, and I get to observe  $Z$ . What have I learned about  $Y$ ?

If  $Y$  and  $Z$  are independent, then the answer is of course nothing. But if they are not independent, then the marginal distribution of  $Y$  changes. In particular, before the random variables were drawn, the (marginal) pdf for  $Y$  was

$$f_Y(y) = \int f_{Y,Z}(y, z) \, dz.$$

After we observe  $Z = z$ , we have

$$f_Y(y|Z = z) = \frac{f_{Y,Z}(y, z)}{f_Z(z)} = \frac{f_{Y,Z}(y, z)}{\int f_{Y,Z}(y, z) \, dy}.$$

$Y$  is still a random variable, but its distribution depends on the value  $z$  that was observed for  $Z$ .

Now, given that I have observed  $Z = z$ , what is the best guess for  $Y$ ? If by “best” we mean that which minimizes the mean squared error, it is the conditional mean. That is, the minimizer of

$$\underset{g}{\text{minimize}} \quad \mathbb{E}[(Y - g)^2 | Z = z]$$

is

$$g = \mathbb{E}[Y | Z = z].$$

Notice that unlike before,  $g$  is not pre-determined, it depends on the outcome  $Z = z$ . We might denote

$$g(z) = E[Y|Z = z].$$

In fact, since  $Z$  is a random variable,  $g$  is a priori also random, we might say

$$g(Z) = E[Y|Z].$$

It is then fair to ask: what is the mean of  $g(Z)$ ? We have<sup>1</sup>

$$\begin{aligned} E[g(Z)] &= E[E[Y|Z]] \\ &= \int E[Y|Z = z] f_Z(z) \, dz \\ &= E[Y] \end{aligned}$$

So on average, we are doing the same thing as if we didn't observe  $Z$  at all, but our MSE will in general be much better.

## Multivariate Gaussian

We say that a random variable  $X \in \mathbb{R}^D$  is a **Gaussian random vector** if there exists a vector  $\boldsymbol{\mu} \in \mathbb{R}^D$  and a symmetric positive definite matrix  $\mathbf{R}$  such that its density can be written as

$$f_X(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} \sqrt{\det(\mathbf{R})}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{R}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right).$$

The vector  $\boldsymbol{\mu}$  is the mean of this distribution, and  $\mathbf{R}$  is the covariance:

$$\boldsymbol{\mu} = E[X], \quad \mathbf{R} = E[(X - \boldsymbol{\mu})(X - \boldsymbol{\mu})^T].$$

---

<sup>1</sup>That  $E[E[Y|Z]] = E[Y]$  is known as the law of *iterated expectation*. The inside  $E$  above is over  $Y$  while the outside one is over  $Z$ .

We will denote this as

$$X \sim \text{Normal}(\boldsymbol{\mu}, \mathbf{R}).$$

The geometry of the density reflects the eigenstructure of  $\mathbf{R}$  — the level-surfaces of the density are ellipsoids with the eigenvectors of  $\mathbf{R}$  as axes, and radii are proportional to the inverses of the eigenvalues.

Example:

[INSERT]

It is a classical fact that  $R[i, j] = 0$  if and only if entries  $X_i$  and  $X_j$  are independent from one another. It is also a fact that if  $\mathbf{A}$  is a  $D \times D$  invertible matrix, then  $\mathbf{A}X$  is also a Gaussian random vector with mean  $\mathbf{A}\boldsymbol{\mu}$  and covariance  $\mathbf{A}\mathbf{R}\mathbf{A}^T$ . There is a classic transformation that makes the components of  $X$  independent from one another. If  $\mathbf{R} = \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^T$  is the eigenvalue decomposition of  $\mathbf{R}$ , then

$$X' = \mathbf{V}^T X$$

is a Gaussian random vector that has independent entries. This follows from the fact that  $\mathbf{V}$  diagonalizes  $\mathbf{R}$ :

$$\mathbf{V}^T \mathbf{R} \mathbf{V} = \boldsymbol{\Lambda}.$$

This transformation is known as the *Karhunen-Loeve* transform.

## Gaussian Estimation

What does observing part of a Gaussian random vector tell us about the part that we do not observe? That is, suppose

$$X \sim \text{Normal}(\mathbf{0}, \mathbf{R}),$$

and then we observe the first  $1, \dots, p$  entries of  $X$  while entries  $p + 1, \dots, D$  stay hidden. We divide  $X$  into

$$X = \begin{bmatrix} X_o \\ X_h \end{bmatrix}, \quad \text{then observe } X_o = \mathbf{x}_o.$$

What is the conditional density for  $X_h | X_o = \mathbf{x}_o$ ?

It turns out that the conditional density is also Gaussian, just with a different mean and different covariance. To see this, we partition the covariance matrix into 4 parts:

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_o & \mathbf{R}_{oh} \\ \mathbf{R}_{oh}^T & \mathbf{R}_h \end{bmatrix}.$$

The upper left corner contains the  $p \times p$  covariance matrix for the random variables that end up being observed, the lower right corner contains the  $D - p \times D - p$  covariance matrix for the unobserved random variables, and  $\mathbf{R}_{oh}$  is the *cross-correlation* matrix, that captures the dependencies between the observed and unobserved random variables.

We can also partition the inverse covariance

$$\mathbf{R}^{-1} = \begin{bmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{D} \end{bmatrix}.$$

Using the Schur complement (see the Technical Details section below), we can write out these blocks of the inverse as

$$\begin{aligned}\mathbf{D} &= (\mathbf{R}_h - \mathbf{R}_{oh}^T \mathbf{R}_o^{-1} \mathbf{R}_{oh})^{-1} \\ \mathbf{C} &= -\mathbf{R}_o^{-1} \mathbf{R}_{oh} \mathbf{D} \\ \mathbf{B} &= (\text{something})\end{aligned}$$

We could write down the expression for  $\mathbf{B}$  if we really wanted to, but it is long and complicated and it ends up that we don't use it. We will use these expressions later, but to ease the notation below, we will stick with  $\mathbf{B}, \mathbf{C}, \mathbf{D}$ .

We can now compute the conditional density using

$$f_{X_h}(\mathbf{x}_h | \mathbf{x}_0) = \frac{f_{X_o, X_h}(\mathbf{x}_o, \mathbf{x}_h)}{f_{X_o}(\mathbf{x}_o)}.$$

The numerator is proportional to

$$\begin{aligned}f_{X_o, X_h}(\mathbf{x}_o, \mathbf{x}_h) &\propto \exp \left( -\frac{1}{2} \begin{bmatrix} \mathbf{x}_o^T & \mathbf{x}_h^T \end{bmatrix} \begin{bmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{x}_o \\ \mathbf{x}_h \end{bmatrix} \right) \\ &= \exp \left( -\frac{1}{2} \left[ \mathbf{x}_o^T \mathbf{B} \mathbf{x}_o + \mathbf{x}_o^T \mathbf{C} \mathbf{x}_h + \mathbf{x}_h^T \mathbf{C}^T \mathbf{x}_o + \mathbf{x}_h^T \mathbf{D} \mathbf{x}_h \right] \right).\end{aligned}$$

The first term above does not depend on  $\mathbf{x}_h$ , so we can write the conditional density as

$$f_{X_h}(\mathbf{x}_h | \mathbf{x}_0) = g(\mathbf{x}_0) \exp \left( -\frac{1}{2} \left[ \mathbf{x}_o^T \mathbf{C} \mathbf{x}_h + \mathbf{x}_h^T \mathbf{C}^T \mathbf{x}_o + \mathbf{x}_h^T \mathbf{D} \mathbf{x}_h \right] \right),$$

where  $g(\mathbf{x}_0)$  is a function that incorporates  $1/f_{X_o}(\mathbf{x}_o)$  and  $\exp(-\mathbf{x}_o^T \mathbf{B} \mathbf{x}_o/2)$  along with some constants. We are not too worried about what  $g$  actually is, just that it does not depend on  $\mathbf{x}_h$ .

To show that  $X_h|X_o$  is a Gaussian random vector, we need a density that looks like

$$(\text{stuff with no } \mathbf{x}_h) \cdot \exp \left( -\frac{1}{2}(\mathbf{x}_h - \boldsymbol{\mu})^T \mathbf{K}(\mathbf{x}_h - \boldsymbol{\mu}) \right).$$

To get our conditional density in this form, we *complete the square* in the exponent. You can easily check that the following relation holds:

$$\begin{aligned} \mathbf{x}_o^T \mathbf{C} \mathbf{x}_h + \mathbf{x}_h^T \mathbf{C}^T \mathbf{x}_o + \mathbf{x}_h^T \mathbf{D} \mathbf{x}_h = \\ (\mathbf{x}_h + \mathbf{D}^{-1} \mathbf{C}^T \mathbf{x}_o)^T \mathbf{D} (\mathbf{x}_h + \mathbf{D}^{-1} \mathbf{C}^T \mathbf{x}_o) - \mathbf{x}_o^T \mathbf{C} \mathbf{D}^{-1} \mathbf{C}^T \mathbf{x}_o. \end{aligned}$$

Again, the last term above does not depend on  $\mathbf{x}_h$ . Thus we can write

$$f_{X_h}(\mathbf{x}_h|\mathbf{x}_o) = h(\mathbf{x}_o) \exp \left( -\frac{1}{2}(\mathbf{x}_h + \mathbf{D}^{-1} \mathbf{C}^T \mathbf{x}_o)^T \mathbf{D} (\mathbf{x}_h + \mathbf{D}^{-1} \mathbf{C}^T \mathbf{x}_o) \right),$$

where  $h(\mathbf{x}_o) = g(\mathbf{x}_o) \exp(\mathbf{x}_o^T \mathbf{C} \mathbf{D}^{-1} \mathbf{C}^T \mathbf{x}_o/2)$ . Plugging in the expressions for  $\mathbf{C}$  and  $\mathbf{D}$  above, we see that  $X_h|X_o = \mathbf{x}_o$  is a Gaussian random vector

$$X_h|X_o = \mathbf{x}_o \sim \text{Normal}(\mathbf{R}_{oh}^T \mathbf{R}_o^{-1} \mathbf{x}_o, \mathbf{R}_h - \mathbf{R}_{oh}^T \mathbf{R}_o^{-1} \mathbf{R}_{oh}).$$

So given the observations  $X_o = \mathbf{x}_o$ , our best (MMSE) guess for  $\mathbf{h}_h$  is the conditional mean:

$$\hat{\mathbf{x}}_h = \mathbf{R}_{oh}^T \mathbf{R}_o^{-1} \mathbf{x}_o.$$

The MSE we will incur with this choice is

$$\mathbb{E}[\|\hat{\mathbf{x}}_h - X_h\|_2^2 | X_o = \mathbf{x}_o] = \text{trace}(\mathbf{R}_h - \mathbf{R}_{oh}^T \mathbf{R}_o^{-1} \mathbf{R}_{oh}).$$

It is a fact that  $\text{trace}(\mathbf{R}_h - \mathbf{R}_{oh}^T \mathbf{R}_o^{-1} \mathbf{R}_{oh}) \leq \text{trace}(\mathbf{R}_h)$  (why?), so the observing  $X_o = \mathbf{x}_o$  also reduces the mean-squared error associated with our best guess.

Notice that for zero-mean Gaussian random variables,  $R[i, j] = 0$  if and only if  $X_i$  and  $X_j$  are independent. Above, this means that if  $X_o$  and  $X_k$  are independent, we will have  $\mathbf{R}_{oh} = \mathbf{0}$ , and the conditional distribution for  $X_h$  is no different from its original marginal (exactly as we would expect).

## Conditional independence

As we just mentioned above, It is easy to interpret a zero-valued entry in the covariance matrix:  $R[i, j] = 0$  means  $X[i]$  and  $X[j]$  are independent. But the entries of the *inverse covariance* matrix also carry interesting information about the dependency structure in  $X$ . Let

$$\mathbf{S} = \mathbf{R}^{-1}.$$

What does it mean if  $S[i, j] = 0$ ? The answer is that  $X[i]$  and  $X[j]$  are independent given observations of all of the other entries  $\{X[k], k \neq i, j\}$  in  $X$ . To see this, suppose that we partition  $\mathbf{S}$  the same way we partitioned  $\mathbf{R}$ :

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_o & \mathbf{S}_{oh} \\ \mathbf{S}_{oh}^T & \mathbf{S}_h \end{bmatrix}$$

We have already seen we can use the *Schur complement* to get an expression for  $\mathbf{S}_h$  in terms of the blocks in  $\mathbf{R}$ :

$$\mathbf{S}_h = (\mathbf{R}_h - \mathbf{R}_{oh}^T \mathbf{R}_o^{-1} \mathbf{R}_{oh})^{-1}.$$

Notice that this is exactly the inverse covariance of the conditional random vector  $X_h|X_o$ . Then, consider the particular case when there are two “hidden” entries in  $X_h$ , say  $X_h = \{X[i], X[j]\}$  and the other  $d - 2$  are in  $X_o$ . The off-diagonal terms in  $\mathbf{S}_h$  above correspond to



$S[i, j]$  and  $S[j, i]$ , if these are zero, then  $\mathbf{S}_h$  is diagonal, and so is  $\mathbf{S}_h^{-1} = \mathbf{R}_{h|o}$ . This means that  $X[i]$  and  $X[j]$  are **conditionally independent** given observations of  $\{X[k], k \neq i, j\}$ .

### Independence and conditional independence

Let  $X \sim \text{Normal}(\boldsymbol{\mu}, \mathbf{R})$ , and set  $\mathbf{S} = \mathbf{R}^{-1}$ . Then

$$R[i, j] = 0 \quad \Leftrightarrow \quad X[i] \text{ and } X[j] \text{ are independent,}$$

and with  $X_{\overline{(i,j)}} = \{X[k], k \neq i, j\}$ ,

$$S[i, j] = 0 \quad \Leftrightarrow \quad X[i]|X_{\overline{(i,j)}} \text{ and } X[j]|X_{\overline{(i,j)}} \text{ are independent.}$$

## Gaussian Graphical Models (a first look)

The non-zero patterns in the inverse covariance  $\mathbf{S}$  often times more descriptive of the dependency structure in  $X$  than the non-zero entries in  $\mathbf{R}$ . For example, let  $Z[1], \dots, Z[d]$  be independent Gaussian random variables with  $Z[k] \sim \text{Normal}(0, 1)$ , and  $0 < a < 1$ , and set

$$\begin{aligned} X[1] &= \sigma Z[1], \\ X[2] &= aX[1] + Z[2] \\ X[3] &= aX[2] + Z[3] \\ &\vdots \\ X[d] &= aX[d-1] + Z[d] \end{aligned}$$

where  $\sigma = (1 - a^2)^{-1/2}$  is chosen so that all the  $X[k]$  have the same variance of  $\sigma^2 = (1 - a^2)^{-1}$ . This is a standard auto-regressive process

— the next point in the vector is computed by taking the previous point, multiplying it by a fixed number, then adding an independent perturbation.

The “flow” of this process is very naturally described with this graph:  
[INSERT]

The covariance matrix for the vector described above is non-zero everywhere:

$$\mathbf{R} = \begin{bmatrix} \sigma^2 & a\sigma^2 & a^2\sigma^2 & \dots & a^{d-1}\sigma^2 \\ a\sigma^2 & \sigma^2 & a\sigma^2 & \dots & a^{d-2}\sigma^2 \\ \vdots & & \ddots & & \vdots \\ a^{d-1}\sigma^2 & a^{d-2}\sigma^2 & \dots & & \sigma^2 \end{bmatrix},$$

but the inverse covariance immediately reveals the structure of the equations that generate  $\mathbf{X}$ :

$$\mathbf{S} = \mathbf{R}^{-1} = \begin{bmatrix} 1 & -a & 0 & \dots & 0 \\ -a & (1+a^2) & -a & \dots & 0 \\ 0 & -a & (1+a^2) & -a & \dots & 0 \\ \vdots & & & & & \\ 0 & \dots & 0 & -a & (1+a^2) & -a \\ 0 & \dots & & 0 & -a & 1 \end{bmatrix}$$

## Technical Details: The Schur Complement

Suppose that  $\mathbf{M}$  is an invertible matrix broken into four blocks:

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_{11} & \mathbf{M}_{12} \\ \mathbf{M}_{21} & \mathbf{M}_{22} \end{bmatrix}.$$

If  $\mathbf{M}_{22}$  is invertible, then the inverse of  $\mathbf{M}$  can be expressed in terms of these blocks using the *Schur complement* of  $\mathbf{M}$  in  $\mathbf{M}_{22}$ :

$$\mathbf{S} = \mathbf{M}_{11} - \mathbf{M}_{12}\mathbf{M}_{22}^{-1}\mathbf{M}_{21}.$$

Then,

$$\mathbf{M}^{-1} = \begin{bmatrix} \mathbf{S}^{-1} & -\mathbf{S}^{-1}\mathbf{M}_{12}\mathbf{M}_{22}^{-1} \\ -\mathbf{M}_{22}^{-1}\mathbf{M}_{21}\mathbf{S}^{-1} & \mathbf{M}_{22}^{-1} + \mathbf{M}_{22}^{-1}\mathbf{M}_{21}\mathbf{S}^{-1}\mathbf{M}_{12}\mathbf{M}_{22}^{-1} \end{bmatrix}$$

Similarly, if  $\mathbf{M}_{11}$  is invertible, we can do something similar with the Schur complement of  $\mathbf{M}$  in  $\mathbf{M}_{11}$ :

$$\mathbf{M}^{-1} = \begin{bmatrix} \mathbf{M}_{11}^{-1} + \mathbf{M}_{11}^{-1}\mathbf{M}_{12}\mathbf{S}^{-1}\mathbf{M}_{21}\mathbf{M}_{11}^{-1} & -\mathbf{M}_{11}^{-1}\mathbf{M}_{12}\mathbf{S}^{-1} \\ -\mathbf{S}^{-1}\mathbf{M}_{21}\mathbf{M}_{11}^{-1} & \mathbf{S}^{-1} \end{bmatrix},$$

where now

$$\mathbf{S} = \mathbf{M}_{22} - \mathbf{M}_{21}\mathbf{M}_{11}^{-1}\mathbf{M}_{12}.$$

These formulas can be checked simply by multiplying out  $\mathbf{M}^{-1}\mathbf{M}$  and seeing that it is  $\mathbf{I}$ .