

III. Statistical Estimation and Classification

In the last chapter of this course, we considered a very specific estimation framework: the unknown variables were put through a linear operator (i.e. a matrix), perturbed, and then observed. In this chapter, we also consider estimating unknown variables from indirect observations, but here our observational model will be completely different. Our observations are samples of a random variable whose distribution is controlled (or parameterized) by the unknown $\boldsymbol{\theta}$:

$$\text{observe } Y_i \sim f_{\boldsymbol{\theta}}(y) \quad \text{for } i = 1, \dots, M.$$

Above, $f_{\boldsymbol{\theta}}$ is the *probability density function* (pdf) for the observations — this density is different for different parameters $\boldsymbol{\theta}$. Moreover, the observations Y_i need not be scalars; they could just as easily be random vectors (or random matrices, or random function, etc). From these observations of Y , we come up with an estimate of $\boldsymbol{\theta}$.

There are two general scenarios in which this kind of problem is treated: we can treat the unknown parameters $\boldsymbol{\theta}$ as being deterministic (they are what they are, we just don't know what they are) or as being random themselves (they obey some probability law that we might know). While there are philosophical differences between these two kinds of treatment, most of the time the application we have dictates which one is more natural.

We start with one example from each of these scenarios; both of these will rely heavily on the linear algebra developed in the last chapter.

Best Linear Unbiased Estimator

Let's see what we get when we treat our canonical linear inverse problem

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{e},$$

as a statistical estimation problem. Our model here is that the error vector \mathbf{e} is random, has zero mean,

$$\mathbb{E}[e[m]] = 0, \quad m = 1, \dots, M, \quad \text{or} \quad \mathbb{E}[\mathbf{e}] = \mathbf{0},$$

and covariance matrix

$$R[\ell, m] = \mathbb{E}[e[\ell]e[m]], \quad \text{or} \quad \mathbb{E}[\mathbf{e}\mathbf{e}^T] = \mathbf{R}.$$

Notice that the diagonal entries of \mathbf{R} contain the variances of the entries of \mathbf{e} , while the off-diagonal terms capture the correlations. We also note here that covariance matrices are always symmetric positive semi-definite, since

$$\mathbf{x}^T \mathbf{R} \mathbf{x} = \mathbf{x}^T \mathbb{E}[\mathbf{e}\mathbf{e}^T] \mathbf{x} = \mathbb{E}[\mathbf{x}^T \mathbf{e}\mathbf{e}^T \mathbf{x}] = \mathbb{E}[|\mathbf{x}^T \mathbf{e}|^2] \geq 0.$$

Since \mathbf{e} is random, our observation \mathbf{y} is random as well. Our unknown parameters \mathbf{x}_0 effect the distribution of \mathbf{y} in an obvious way: \mathbf{y} has the same distribution as \mathbf{e} , only the mean is $\mathbf{A}\mathbf{x}_0$ instead of $\mathbf{0}$.

We can now ask what is the best way to estimate \mathbf{x}_0 . There are all kinds of ways in which we might treat this problem, but for now we will restrict ourselves to estimators that have the following properties:

1. **Linearity.** That is, our estimate can be computed by applying a fixed matrix to \mathbf{y} ,

$$\hat{\mathbf{x}} = \mathbf{L}\mathbf{y},$$

for some $N \times M$ matrix \mathbf{L} .

2. **Unbiased.** Since the estimate $\hat{\mathbf{x}}$ is a function of random variables, it is itself a random variable. Our estimator is unbiased if

$$\mathbb{E}[\hat{\mathbf{x}}] = \mathbf{x}_0,$$

which means the expectation of the estimation error is zero,

$$\mathbb{E}[\hat{\mathbf{x}} - \mathbf{x}_0] = \mathbf{0}.$$

We will search for the best such estimator; the best linear unbiased estimator (BLUE). And by “best” we mean that the mean-squared estimation error (MSE), $\mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2]$ is minimized.

Our derivation of the BLUE will require some manipulation of covariance matrices, so let’s do some warm-up exercises first.

Questions:

1. Suppose that the entries of \mathbf{e} have variances $\nu_m^2 = \mathbb{E}[e[m]^2]$. Calculate

$$\mathbb{E}[\|\mathbf{e}\|_2^2] = \underline{\hspace{2cm}}.$$

(the expected energy of \mathbf{e}).

Answer:

$$\begin{aligned}\mathbb{E}[\|\mathbf{e}\|_2^2] &= \sum_{m=1}^M \mathbb{E}[e[m]^2] \\ &= \sum_{m=1}^M \nu_m^2.\end{aligned}$$

2. Now let \mathbf{D} be a diagonal matrix

$$\mathbf{D} = \begin{bmatrix} d_1 & & & \\ & d_2 & & \\ & & \ddots & \\ & & & d_M \end{bmatrix}.$$

Calculate

$$E[\|\mathbf{D}\mathbf{e}\|_2^2] = \underline{\hspace{2cm}}.$$

Answer:

$$\begin{aligned} E[\|\mathbf{D}\mathbf{e}\|_2^2] &= \sum_{m=1}^M E[d_m^2 e[m]^2] \\ &= \sum_{m=1}^M d_m^2 \nu_m^2. \end{aligned}$$

3. Suppose that $\mathbf{e} \in \mathbb{R}^M$ is a random vector with zero mean and covariance matrix \mathbf{R} . Let \mathbf{M} be an arbitrary $N \times M$ matrix, and set $\mathbf{g} = \mathbf{M}\mathbf{e}$. Clearly $E[\mathbf{g}] = \mathbf{0}$. What is the covariance $E[\mathbf{g}\mathbf{g}^T]$?

Answer:

$$\begin{aligned} E[\mathbf{g}\mathbf{g}^T] &= E[\mathbf{M}\mathbf{e}(\mathbf{M}\mathbf{e})^T] \\ &= E[\mathbf{M}\mathbf{e}\mathbf{e}^T\mathbf{M}^T] = \mathbf{M} E[\mathbf{e}\mathbf{e}^T] \mathbf{M}^T = \mathbf{M}\mathbf{R}\mathbf{M}^T. \end{aligned}$$

4. Suppose $\mathbf{e} \in \mathbb{R}^M$ has covariance matrix \mathbf{R} . Let \mathbf{L} be an $N \times M$ matrix. Calculate

$$E[\|\mathbf{L}\mathbf{e}\|_2^2] = \underline{\hspace{2cm}}.$$

Answer: We use two facts which are easily verified (do this at home). First, the inner product of two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^N$ is equal to the trace of their outer product:

$$\mathbf{v}^T \mathbf{u} = \text{trace}(\mathbf{u}\mathbf{v}^T).$$

Second, if \mathbf{Q} is a square matrix whose entries are random variables, then

$$\mathbb{E}[\text{trace}(\mathbf{Q})] = \text{trace}(\mathbb{E}[\mathbf{Q}]).$$

Then

$$\begin{aligned} \mathbb{E}[\|\mathbf{L}\mathbf{e}\|_2^2] &= \mathbb{E}[\langle \mathbf{L}\mathbf{e}, \mathbf{L}\mathbf{e} \rangle] \\ &= \mathbb{E}[\text{trace}(\mathbf{L}\mathbf{e}\mathbf{e}^T \mathbf{L}^T)] \\ &= \text{trace} \left(\mathbb{E}[\mathbf{L}\mathbf{e}\mathbf{e}^T \mathbf{L}^T] \right) \\ &= \text{trace} \left(\mathbf{L} \mathbb{E}[\mathbf{e}\mathbf{e}^T] \mathbf{L}^T \right) \\ &= \text{trace}(\mathbf{L}\mathbf{R}\mathbf{L}^T). \end{aligned}$$

We return now to the derivation of the BLUE for \mathbf{x}_0 from observations $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{e}$. We are restricting ourselves to linear estimators, so we can write

$$\hat{\mathbf{x}} = \mathbf{L}\mathbf{y} = \mathbf{L}(\mathbf{A}\mathbf{x} + \mathbf{e}) = \mathbf{L}\mathbf{A}\mathbf{x} + \mathbf{L}\mathbf{e},$$

for some matrix \mathbf{L} which we will optimize. We want the estimator to be unbiased, so

$$\begin{aligned} \mathbf{0} &= \mathbb{E}[\mathbf{x}_0 - \hat{\mathbf{x}}] = \mathbb{E}[\mathbf{x}_0 - \mathbf{L}\mathbf{A}\mathbf{x} - \mathbf{L}\mathbf{e}] \\ &= \mathbf{x}_0 - \mathbf{L}\mathbf{A}\mathbf{x}_0 - \mathbb{E}[\mathbf{L}\mathbf{e}] \\ &= \mathbf{x}_0 - \mathbf{L}\mathbf{A}\mathbf{x}_0, \end{aligned}$$

where the last step comes from the fact that $\mathbb{E}[\mathbf{L}\mathbf{e}] = \mathbf{0}$, since $\mathbb{E}[\mathbf{e}] = \mathbf{0}$. Thus we need \mathbf{L} to obey

$$\mathbf{L}\mathbf{A}\mathbf{x}_0 = \mathbf{x}_0.$$

That is, we want \mathbf{L} to be a **left inverse** of \mathbf{A} , meaning $\mathbf{L}\mathbf{A} = \mathbf{I}$.

With these two properties in hand, the variance of our estimate for a qualifying \mathbf{L} is

$$\begin{aligned} \mathbb{E}[\|\mathbf{x}_0 - \hat{\mathbf{x}}\|_2^2] &= \mathbb{E}[\|\mathbf{x}_0 - \mathbf{L}\mathbf{A}\mathbf{x}_0 - \mathbf{L}\mathbf{e}\|_2^2] \\ &= \mathbb{E}[\|\mathbf{L}\mathbf{e}\|_2^2] \\ &= \mathbb{E}[\text{trace}(\mathbf{L}\mathbf{R}\mathbf{L}^T)]. \end{aligned}$$

So we would like to find the matrix which minimizes

$$\underset{\mathbf{L} \in \mathbb{R}^{N \times M}}{\text{minimize}} \quad \text{trace}(\mathbf{L}\mathbf{R}\mathbf{L}^T) \quad \text{subject to} \quad \mathbf{L}\mathbf{A} = \mathbf{I}.$$

I propose that the solution to the above is

$$\mathbf{L}_0 = (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1}.$$

Let's check this. Clearly $\mathbf{L}_0 \mathbf{A} = \mathbf{I}$, so \mathbf{L}_0 is a left inverse. It remains to show that for any other left inverse \mathbf{L} ,

$$\text{trace}(\mathbf{L}\mathbf{R}\mathbf{L}^T) \geq \text{trace}(\mathbf{L}_0 \mathbf{R} \mathbf{L}_0^T).$$

Write a candidate \mathbf{L} as

$$\mathbf{L} = \mathbf{L}_0 + (\mathbf{L} - \mathbf{L}_0).$$

Then

$$\begin{aligned} \text{trace}(\mathbf{L}\mathbf{R}\mathbf{L}^T) &= \text{trace}(\mathbf{L}_0 \mathbf{R} \mathbf{L}_0^T) + \text{trace}(\mathbf{L}_0 \mathbf{R} (\mathbf{L} - \mathbf{L}_0)^T) \\ &\quad + \text{trace}((\mathbf{L} - \mathbf{L}_0) \mathbf{R} \mathbf{L}_0^T) + \text{trace}((\mathbf{L} - \mathbf{L}_0) \mathbf{R} (\mathbf{L} - \mathbf{L}_0)^T). \end{aligned}$$

Note that

$$\mathbf{R} \mathbf{L}_0^T = \mathbf{R} \mathbf{R}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} = \mathbf{A} (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1}.$$

Thus

$$(\mathbf{L} - \mathbf{L}_0)\mathbf{R}\mathbf{L}_0^T = (\mathbf{L} - \mathbf{L}_0)\mathbf{A}(\mathbf{A}^T\mathbf{R}^{-1}\mathbf{A})^{-1} = \mathbf{0}$$

since both $\mathbf{L}\mathbf{A} = \mathbf{I}$ and $\mathbf{L}_0\mathbf{A} = \mathbf{I}$. We are left with

$$\text{trace}(\mathbf{L}\mathbf{R}\mathbf{L}^T) = \text{trace}(\mathbf{L}_0\mathbf{R}\mathbf{L}_0^T) + \text{trace}((\mathbf{L} - \mathbf{L}_0)\mathbf{R}(\mathbf{L} - \mathbf{L}_0)^T).$$

Since $(\mathbf{L} - \mathbf{L}_0)\mathbf{R}(\mathbf{L} - \mathbf{L}_0)^T$ is symmetric and positive semi-definite, the term on the right is ≥ 0 . So we conclude

$$\text{trace}(\mathbf{L}\mathbf{R}\mathbf{L}^T) \geq \text{trace}(\mathbf{L}_0\mathbf{R}\mathbf{L}_0^T) \quad \text{for all left inverses } \mathbf{L}.$$

Best Linear Unbiased Estimator (BLUE):

From observations,

$$\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{e}, \quad \mathbb{E}[\mathbf{e}\mathbf{e}^T] = \mathbf{R},$$

the BLUE is

$$\hat{\mathbf{x}}_{\text{blue}} = (\mathbf{A}^T\mathbf{R}^{-1}\mathbf{A})^{-1}\mathbf{A}^T\mathbf{R}^{-1}\mathbf{y}.$$

A quick calculation shows

$$\mathbf{L}_0\mathbf{R}\mathbf{L}_0^T = (\mathbf{A}^T\mathbf{R}^{-1}\mathbf{A})^{-1},$$

and so the MSE of the BLUE is

$$\begin{aligned} \mathbb{E}[\|\mathbf{x}_0 - \hat{\mathbf{x}}_{\text{blue}}\|_2^2] &= \text{trace}((\mathbf{A}^T\mathbf{R}^{-1}\mathbf{A})^{-1}) \\ &= \text{sum of eigenvalues of } (\mathbf{A}^T\mathbf{R}^{-1}\mathbf{A})^{-1}. \end{aligned}$$

Uncorrelated errors

Suppose that the random errors are uncorrelated, so that the covariance matrix is diagonal

$$\mathbf{R} = \mathbb{E}[\mathbf{e}\mathbf{e}^T] = \begin{bmatrix} \nu_1^2 & 0 & 0 & \cdots \\ 0 & \nu_2^2 & 0 & \cdots \\ \vdots & & \ddots & \\ & & & \nu_M^2 \end{bmatrix}$$

If ν_m is large, it means that we do not have much confidence in our measurement y_m . On the other hand, if ν_m is small, it means that our measurement y_m is most likely very close to the true value of $(\mathbf{A}\mathbf{x}_0)[m]$

The BLUE we derived above tells us exactly how we should weight these measurements. The BLUE is

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1} \mathbf{y},$$

where

$$\mathbf{R}^{-1} = \begin{bmatrix} \frac{1}{\nu_1^2} & 0 & 0 & \cdots \\ 0 & \frac{1}{\nu_2^2} & 0 & \cdots \\ \vdots & & \ddots & \\ & & & \frac{1}{\nu_M^2} \end{bmatrix}.$$

Thus we can interpret the BLUE as the solution to the least-squares problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{A}\mathbf{x}\|_2^2,$$

where \mathbf{W} is a diagonal weighting matrix $W[m, m] = 1/\nu_m$. This is penalizing mismatch in the components in which we have confidence (ν_m small) more than the components in which we do not have confidence (ν_m larger).

Example. We take M readings of a patient's pulse, each has an error of ν^2 . In this case, the underlying quantity (the pulse) x_0 is a scalar. The optimal estimate (no matter what ν is) is

$$\hat{x} = \frac{1}{M} (y[1] + y[2] + \cdots + y[M]) .$$

What is the mean-square error for this estimate?

Answer: The mean-square error is

$$\begin{aligned} \mathbb{E}[|x_0 - \hat{x}|^2] &= \mathbb{E} \left[\left| x_0 - \frac{1}{M} \sum_{m=1}^M (x_0 + e[m]) \right|^2 \right] \\ &= \mathbb{E} \left[\left| \frac{1}{M} \sum_{m=1}^M e[m] \right|^2 \right] \\ &= \frac{1}{M^2} \mathbb{E}[\langle \mathbf{e}, \mathbf{e} \rangle] \\ &= \frac{1}{M^2} \mathbb{E}[\text{trace}(\mathbf{e}\mathbf{e}^T)] \\ &= \frac{1}{M^2} \text{trace}(\mathbb{E}[\mathbf{e}\mathbf{e}^T]) \\ &= \frac{\nu^2}{M}, \end{aligned}$$

where the last step follows from the fact that the covariance matrix of the errors \mathbf{e} is diagonal.

Now suppose that the variance for each of the M measurements is different; $\nu_1^2, \nu_2^2, \dots, \nu_M^2$.

Now what is the best estimate \hat{x} ?

What is the MSE of this estimate?

Answers: We have

$$\mathbf{A} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \quad \mathbf{R}^{-1} = \begin{bmatrix} 1/\nu_1^2 & & & \\ & 1/\nu_2^2 & & \\ & & \ddots & \\ & & & 1/\nu_M^2 \end{bmatrix},$$

and

$$(\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} = \left(\sum_{m=1}^M 1/\nu_m^2 \right)^{-1},$$

and so

$$\hat{x} = \frac{\sum_{m=1}^M y[m]/\nu_m^2}{\sum_{m=1}^M 1/\nu_m^2}.$$

The MSE is

$$\text{trace}((\mathbf{A} \mathbf{R}^{-1} \mathbf{A})^{-1}) = \left(\sum_{m=1}^M 1/\nu_m^2 \right)^{-1}.$$

Exercise with correlated errors: We measure

$$\mathbf{y} = \mathbf{A} \mathbf{x} + \mathbf{e}$$

with

$$\mathbf{A} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 0 \\ 5 \end{bmatrix}, \quad \text{E}[\mathbf{e} \mathbf{e}^T] = \mathbf{R} = \begin{bmatrix} 3 & -1 \\ -1 & 2 \end{bmatrix}.$$

1. Find the best linear unbiased estimate.

Hint:

$$\mathbf{R}^{-1} = \frac{1}{5} \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix}.$$

2. Calculate $E[\|\mathbf{x}_0 - \hat{\mathbf{x}}_{\text{blue}}\|_2^2]$.

Three closing notes on the BLUE

1. Notice that we did not need to model the entire probability distribution of \mathbf{e} , we only needed knowledge of its first two moments (mean and covariance). We got away with this because we restricted our estimator to be linear and unbiased — other classes of estimators will in general rely on higher-order moments or even knowledge of the complete distribution of \mathbf{e} .
2. The assumption that the estimator is unbiased, that $E[\hat{\mathbf{x}}] = \mathbf{x}_0$, seemed natural. However, there can be significant gains (in terms of MSE) from loosening this restriction — introducing a little bias can allow us to reduce the overall error.
3. The BLUE is overall optimal (in terms of MSE) when \mathbf{e} is a Gaussian random vector. That is, even if we allow ourselves to look at nonlinear, biased estimators, we will not find one with better performance. It is no coincidence that the probability law for Gaussian random vectors is completely determined by a mean and a covariance matrix.