

Weighted Least Squares

Standard least-squares tries to fit a vector \mathbf{x} to a set of “measurements” \mathbf{y} by solving

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

Now, what if some of the measurements are more reliable than others? Or, what if the errors are closely correlated between measurements?

There is a systematic way to treat both of these cases using **weighted least-squares**. Instead of minimizing the energy in the residual

$$\|\mathbf{r}\|_2^2 = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2,$$

we will minimize

$$\|\mathbf{W}\mathbf{r}\|_2^2 = \|\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{A}\mathbf{x}\|_2^2,$$

for some $M \times M$ **weighting matrix** \mathbf{W} .

When \mathbf{W} is a diagonal matrix,

$$\mathbf{W} = \begin{bmatrix} w_{11} & & & \\ & w_{22} & & \\ & & \ddots & \\ & & & w_{MM} \end{bmatrix},$$

then the error we are minimizing looks like

$$\|\mathbf{W}\mathbf{r}\|_2^2 = w_{11}^2 r[1]^2 + w_{22}^2 r[2]^2 + \cdots + w_{MM}^2 r[M]^2.$$

By adjusting the w_{mm} , we can penalize some of the components of the error more than others.

By adding off-diagonal terms, we can account for **correlations** in the error (we will explore this further later in these notes).

Solving

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

is simple. We simply use least-squares with $\mathbf{W}\mathbf{A}$ as the matrix, and $\mathbf{W}\mathbf{y}$ as the observations:

$$\hat{\mathbf{x}}_{\text{wls}} = (\mathbf{W}\mathbf{A})^\dagger \mathbf{W}\mathbf{y},$$

where $(\mathbf{W}\mathbf{A})^\dagger$ is the pseudo-inverse of $\mathbf{W}\mathbf{A}$.

For the rest of this section, we will assume that $M \geq N$ (meaning that there are at least as many observations as unknowns) and that \mathbf{A} has full column rank. This allows us to write

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

Example: We measure a patient's pulse 3 times, and record

$$y[1] = 70, \quad y[2] = 80, \quad y[3] = 120.$$

In this case, we can take

$$\mathbf{A} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

What is the least-square estimate for the pulse rate x_0 ?

Now say that we were in a hurry when the third measurement was made, so we would like to weigh less than the others. What is the weighted least-squares estimate when

$$\mathbf{W} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & w_{33} \end{bmatrix} ?$$

What about the particular case when $w_{33} = 1/2$?

Statistical Estimation

Suppose we use the following model for our measurements:

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

where $\mathbf{y} \in \mathbb{R}^M$, \mathbf{A} is an $M \times N$ matrix, $\mathbf{x}_0 \in \mathbb{R}^N$ is what we are interested in estimating, and $\mathbf{e} \in \mathbb{R}^M$ is a **random error**.

We will assume that each entry of \mathbf{e} has zero mean:

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

We will characterize \mathbf{e} through its **covariance matrix**

$$R[\ell, m] = \mathbb{E}[e[\ell]e[m]],$$

or more compactly

$$\mathbf{R} = \mathbb{E}[\mathbf{e}\mathbf{e}^T].$$

The diagonal of \mathbf{R} contains the variances of the entries of \mathbf{e} , while the off diagonal terms capture the correlations (which is the same as covariance, since all of the $e[m]$ are zero mean).

For example, if two measurement errors have

$$\begin{aligned} \text{var}(e[1]) = \mathbb{E}[e[1]^2] &= 3, & \text{var}(e[2]) = \mathbb{E}[e[2]^2] &= 2, \\ \text{and } \text{cov}(e[1], e[2]) &= \mathbb{E}[e[1]e[2]] = -1, \end{aligned}$$

then

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

It is always true that covariance matrices are symmetric and positive semi-definite (so their eigenvalues are ≥ 0).

A handy fact that we will use repeatedly below is that if \mathbf{e} has covariance matrix \mathbf{R} , then for any matrix \mathbf{M} , the covariance of $\mathbf{M}\mathbf{e}$ is¹

$$\mathbb{E}[\mathbf{M}\mathbf{e}(\mathbf{M}\mathbf{e})^T] = \mathbb{E}[\mathbf{M}\mathbf{e}\mathbf{e}^T\mathbf{M}^T] = \mathbf{M}\mathbb{E}[\mathbf{e}\mathbf{e}^T]\mathbf{M}^T = \mathbf{M}\mathbf{R}\mathbf{M}^T.$$

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:

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

¹If you want to see why that second-to-last step is true more explicitly, set $\mathbf{Q} = \mathbf{M}\mathbf{e}\mathbf{e}^T\mathbf{M}^T$. Then if \mathbf{m}_i is the i th row of \mathbf{M} ,

$$Q[i, j] = (\mathbf{M}\mathbf{e})[i](\mathbf{M}\mathbf{e})[j] = \langle \mathbf{e}, \mathbf{m}_i \rangle \langle \mathbf{m}_j, \mathbf{e} \rangle = \sum_{\ell} \sum_k M[i, \ell] M[j, k] e[\ell] e[k],$$

and

$$\mathbb{E}[Q[i, j]] = \sum_{\ell} \sum_k M[i, \ell] M[j, k] R[\ell, k] = (\mathbf{M}\mathbf{R}\mathbf{M}^T)[i, j],$$

so $\mathbb{E}[\mathbf{Q}] = \mathbf{M}\mathbf{R}\mathbf{M}^T$.

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

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

Calculate

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

Answer:

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

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

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

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:

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

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

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

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 $(\mathbf{A}\mathbf{x}_0)[m]$

We will see this rigorously below, but in this case, the “correct” weighting for each component is simply the inverse of the standard deviation; the weighting matrix \mathbf{W} should be diagonal with

$$W[m, m] = \frac{1}{\nu_m}, \quad (\mathbf{W} = \mathbf{R}^{-1/2}).$$

Then the weighted least-squares estimate is given by

$$\begin{aligned}\hat{\mathbf{x}}_{\text{wls}} &= (\mathbf{A}^T \mathbf{W}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W}^T \mathbf{W} \mathbf{y} \\ &= (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1} \mathbf{y}.\end{aligned}$$

The reconstruction error of this estimate is

$$\begin{aligned}\mathbf{x}_0 - \hat{\mathbf{x}}_{\text{wls}} &= \mathbf{x}_0 - (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1} (\mathbf{A} \mathbf{x}_0 + \mathbf{e}) \\ &= -(\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1} \mathbf{e}\end{aligned}$$

The **mean-square error** (MSE) of the error for this estimate is calculated using the result of Question 3 above:

$$\begin{aligned}\mathbb{E}[\|\mathbf{x}_0 - \hat{\mathbf{x}}_{\text{wls}}\|_2^2] &= \text{trace} \left((\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1} \mathbf{R} \mathbf{R}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \right) \\ &= \text{trace} \left((\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{R}^{-1} \mathbf{A} (\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \right) \\ &= \text{trace} \left((\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1} \right)\end{aligned}$$

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

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

Best Linear Unbiased Estimator (BLUE)

We now return to the general estimation problem: we observe

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

where $\mathbf{e} \in \mathbb{R}^M$ is random with

$$\mathbb{E}[\mathbf{e}] = \mathbf{0}, \quad \mathbb{E}[\mathbf{e}\mathbf{e}^T] = \mathbf{R}.$$

Since \mathbf{e} is random, the observations \mathbf{y} are also random. We can now ask what is the best statistical estimate of \mathbf{x}_0 . 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).

Let's make it clear what we mean by "best". We mean that the MSE of the estimation error, $\mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2]$ is minimized.

The estimator is linear, 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} ,

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

Write a candidate \mathbf{L} as

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

Then

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

Note that

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

Thus

$$\begin{aligned} (\mathbf{L} - \mathbf{L}_0)\mathbf{RL}_0^T &= (\mathbf{L} - \mathbf{L}_0)\mathbf{A}(\mathbf{A}^T\mathbf{R}^{-1}\mathbf{A})^{-1} \\ &= \mathbf{0} \end{aligned}$$

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

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

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

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

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

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

$(\mathbf{A}^T \mathbf{R}^{-1} \mathbf{A})^{-1}$ is sometimes called the **information matrix**.

Exercise: We measure

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

with

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

1. Find the best linear unbiased estimate.

Hint:

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

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