

X is a random variable whose value we want to estimate. We generally have information available on which to base our estimate, typically other random variables whose values we observe. Thus our estimate is in general a function of other random variables and is therefore itself a random variable. In constructing our estimate we might be allowed to use only certain functions, e.g. linear or affine, of the observed random variables. In any event, if we denote our estimate \hat{X} , the mean squared error of our estimate is

$$(1) \quad \mathbb{E} \left((X - \hat{X})^2 \right),$$

and the goal of any minimum mean squared estimation error (MMSE) estimation scheme is to find \hat{X} that minimizes this quantity over all estimates that satisfy our information and functional-implementation constraints. In what follows, I'll consider a few different sets of constraints and see how they affect the solution to the problem and the resulting minimum value of (1).

We can always decompose the mean squared error (1) as

$$(2) \quad \mathbb{E} \left((X - \hat{X})^2 \right) = \text{Var}(X - \hat{X}) + (\mathbb{E}(X - \hat{X}))^2.$$

The first term is the variance of the estimation error $\tilde{X} = X - \hat{X}$ and the second is the square of the expected value of \tilde{X} . That expected value is known as the bias of the estimate \hat{X} , and we call \hat{X} an unbiased estimate when $\mathbb{E}(\tilde{X}) = 0$. All the MMSE estimators treated in this handout are unbiased and at the same time minimize the variance of the estimation error over all estimators satisfying the information and implementation constraints. It makes sense that a good MMSE estimator will have zero expected estimation error and small estimation-error variance, which together imply that the estimation error's probability mass is concentrated heavily around zero, but implementation constraints sometimes preclude simultaneous minimization of bias and estimation-error variance.

The most restrictive information constraint arises in the trivial situation where we have no extra information at all. In that case, our estimate \hat{X} must be a constant. Our MMSE estimation problem demands that we choose a number \hat{b} to minimize

$$\mathbb{E} \left((X - b)^2 \right)$$

over all real numbers b . Perhaps it's obvious to you that the optimal choice of b is $\hat{b} = \mathbb{E}(X)$, and indeed that's the case, but let's prove it. For any choice of b , equation (2) reads

$$\begin{aligned} \mathbb{E} \left((X - b)^2 \right) &= \text{Var}(X - b) + (\mathbb{E}(X - b))^2 \\ &= \text{Var}(X) + (\mathbb{E}(X) - b)^2, \end{aligned}$$

Thus our mean squared error is the sum of $\text{Var}(X)$, a term we can do nothing about, and $(\mathbb{E}(X) - b)^2$, which we can zero out by setting $b = \mathbb{E}(X)$. Summarize as follows.

Constant MMSE Estimation: The constant \hat{X}_{opt} that minimizes (1) over all constants \hat{X} is $\hat{X}_{\text{opt}} = \mathbb{E}(X)$.

Now suppose we have access to one other random variable Y , and we're allowed to choose as our estimate $\hat{X} = h(Y)$, where $h(Y)$ is any function of Y . In class, we showed that the minimizing choice of \hat{X} in (1) is $\hat{X}_{\text{opt}} = \mathbb{E}(X | Y)$. Recall how the argument

went. For any function $h(Y)$, we have

$$\begin{aligned}\mathbb{E}((X - h(Y))^2) &= \mathbb{E}((X - \mathbb{E}(X | Y) + \mathbb{E}(X | Y) - h(Y))^2) \\ &= \mathbb{E}((X - \mathbb{E}(X | Y))^2) + 2\mathbb{E}((X - \mathbb{E}(X | Y))(\mathbb{E}(X | Y) - h(Y))) \\ &\quad + \mathbb{E}((\mathbb{E}(X | Y) - h(Y))^2) .\end{aligned}$$

The middle term is zero because, by the law of iterated expectations,

$$\begin{aligned}\mathbb{E}((X - \mathbb{E}(X | Y))(\mathbb{E}(X | Y) - h(Y))) &= \mathbb{E}(\mathbb{E}(X - \mathbb{E}(X | Y))(\mathbb{E}(X | Y) - h(Y)) | Y) \\ &= \mathbb{E}((\mathbb{E}(X | Y) - h(Y))\mathbb{E}(X - \mathbb{E}(X | Y) | Y)) \\ &= \mathbb{E}(\mathbb{E}(X | Y) - h(Y)) \times 0 \\ &= 0 ,\end{aligned}$$

where the second line holds because $\mathbb{E}(X | Y) - h(Y)$ is a function of Y and the third because $\mathbb{E}(\mathbb{E}(X | Y) | Y) = \mathbb{E}(X | Y)$. Since the middle term is zero, we have

$$(3) \quad \mathbb{E}((X - h(Y))^2) = \mathbb{E}((X - \mathbb{E}(X | Y))^2) + \mathbb{E}((\mathbb{E}(X | Y) - h(Y))^2) .$$

This makes it obvious that we must choose $h(Y) = \mathbb{E}(X | Y)$ to minimize the mean squared error. That choice zeroes the second term in (3) and makes the mean squared error equal to $\mathbb{E}((X - \mathbb{E}(X | Y))^2)$, a term we can't affect by choice of h . Here's the summary of this important result.

General MMSE Estimation: The random variable \hat{X}_{opt} that minimizes (1) over all functions of Y is $\hat{X}_{\text{opt}} = \mathbb{E}(X | Y)$.

One noteworthy feature of the general MMSE is that it is unbiased, i.e.

$$\mathbb{E}(\hat{X}_{\text{opt}}) = \mathbb{E}(\mathbb{E}(X | Y)) = \mathbb{E}(X) ,$$

which is equivalent to

$$\mathbb{E}(\tilde{X}) = 0 ,$$

where $\tilde{X} = X - \hat{X}_{\text{opt}}$ is the optimal estimation error. Furthermore, the optimal estimation error is uncorrelated with Y and in fact with any function $h(Y)$. That's because

$$\begin{aligned}\text{Cov}(\tilde{X}_{\text{opt}}, h(Y)) &= \mathbb{E}(\tilde{X}_{\text{opt}}(h(Y) - \tilde{X}_{\text{opt}})) \\ &= \mathbb{E}(\mathbb{E}(\tilde{X}_{\text{opt}}(h(Y) - \tilde{X}_{\text{opt}}) | Y)) \\ &= \mathbb{E}((h(Y) - \tilde{X}_{\text{opt}})\mathbb{E}(\tilde{X}_{\text{opt}} | Y)) \\ &= 0 ,\end{aligned}$$

where the first line holds because \tilde{X} has zero mean, the second by the law of iterated expectations, the third because $h(Y) - \mathbb{E}(h(Y))$ is a function of Y , and the fourth because the inner conditional expectation on the third line is zero.

An identical argument leads to an extended General MMSE Estimation result that applies when multiple random variables are available to estimate X in terms of.

Extended General MMSE Estimation: The random variable \hat{X}_{opt} that minimizes (1) over all functions of Y_1, Y_2, \dots, Y_n is $\hat{X}_{\text{opt}} = \mathbb{E}(X | Y_1, Y_2, \dots, Y_n)$.

We find in this case that the estimation error $\tilde{X} = X - \hat{X}_{\text{opt}}$ has zero mean as before and is uncorrelated with any function $h(Y_1, Y_2, \dots, Y_n)$.

Suppose now that we want to estimate X in terms of Y using a linear — strictly speaking, affine — estimator of the form

$$\hat{X} = aY + b.$$

The mean square error for such an estimator, using (2), is

$$\mathbb{E}((X - aY - b)^2) = \text{Var}(X - aY - b) + (\mathbb{E}(X - aY - b))^2.$$

Regarding the variance,

$$\begin{aligned} \text{Var}(X - aY - b) &= \text{Var}(X - aY) \\ &= \text{Var}(X) - 2a\text{Cov}(X, Y) + a^2\text{Var}(Y), \end{aligned}$$

The value of a minimizing this quantity, which you can calculate by taking the derivative, is

$$\hat{a} = \frac{\text{Cov}(X, Y)}{\text{Var}(Y)},$$

and by choosing

$$b = \hat{b} = \mathbb{E}(X) - \hat{a}\mathbb{E}(Y)$$

we can zero out the bias term in the mean square error expression. Summarize as follows.

Linear MMSE Estimation I: The random variable \hat{X}_{opt} that minimizes (1) over all linear estimators $\hat{X} = aY + b$ is

$$\begin{aligned} \hat{X}_{\text{opt}} &= \hat{a}Y + \hat{b} \\ &= \frac{\text{Cov}(X, Y)}{\text{Var}(Y)}Y + \left(\mathbb{E}(X) - \frac{\text{Cov}(X, Y)}{\text{Var}(Y)}\mathbb{E}(Y) \right) \\ &= \frac{\text{Cov}(X, Y)}{\text{Var}(Y)}(Y - \mathbb{E}(Y)) + \mathbb{E}(X). \end{aligned}$$

Again in this case the optimal estimator is unbiased in the sense that the expected value of \hat{X}_{opt} is the same as the expected value of X . Furthermore, the optimal estimation error $\tilde{X} = X - \hat{X}_{\text{opt}}$ is again uncorrelated with Y , and in fact with any linear function $aY + b$. That's because for any a and b we have

$$\begin{aligned} \text{Cov}(\tilde{X}, aY + b) &= \text{Cov}(X - \hat{a}Y - \hat{b}, aY + b) \\ &= a \text{Cov}(X - \hat{a}Y, Y) \\ &= a (\text{Cov}(X, Y) - \hat{a}\text{Var}(Y)) \\ &= 0, \end{aligned}$$

where the last line holds because $\hat{a} = \text{Cov}(X, Y)/\text{Var}(Y)$.

Suppose now that we want to estimate X linearly in terms of several random variables. First let's consider the case where W_1, \dots, W_n are uncorrelated random variables and we want to minimize (1) using a linear estimator of the form

$$\hat{X} = \sum_{k=1}^n a_k W_k + b.$$

Conveniently, it turns out that the optimal coefficients are

$$\hat{a}_k = \frac{\text{Cov}(X, W_k)}{\text{Var}(W_k)} \quad 1 \leq k \leq n$$

and

$$\hat{b} = \mathbb{E}(X) = \sum_{k=1}^n \hat{a}_k \mathbb{E}(W_k) .$$

Thus the optimal estimator is

$$\begin{aligned} \hat{X}_{\text{opt}} &= \sum_{k=1}^n \hat{a}_k W_k + \hat{b} \\ &= \left(\sum_{k=1}^n \frac{\text{Cov}(X, W_k)}{\text{Var}(W_k)} (W_k - \mathbb{E}(W_k)) \right) + \mathbb{E}(X) . \end{aligned}$$

You can prove this in any number of ways, for example by induction, but here's a simple way to verify that these coefficients do the job. For any a_k , $1 \leq k \leq n$, and b , we have by equation (2)

$$\mathbb{E} \left(\left(X - \sum_{k=1}^n a_k W_k - b \right)^2 \right) = \text{Var} \left(X - \sum_{k=1}^n a_k W_k - b \right) + \left(\mathbb{E} \left(X - \sum_{k=1}^n a_k W_k - b \right) \right)^2 .$$

Take this decomposition of the mean squared error and proceed as follows.

- Rewrite the variance term as

$$\text{Var} \left(X - Q + Q - \sum_{k=1}^n a_k W_k \right) ,$$

where Q is shorthand for \hat{X}_{opt} above. The b disappears because it doesn't affect the variance.

- Note that $\mathbb{E}(Q) = \mathbb{E}(X)$, so $\mathbb{E}(X - Q) = 0$.
- Prove that $X - Q$ is uncorrelated with W_m for $1 \leq m \leq n$. This is because

$$\begin{aligned} \text{Cov}(X - Q, W_m) &= \text{Cov}((X - \hat{a}_m W_m), W_m) \\ &= \text{Cov}(X, W_m) - \hat{a}_m \text{Var}(W_m) \\ &= 0 , \end{aligned}$$

where the first line holds because the W_k are uncorrelated, so the only term in Q not uncorrelated with W_m is the $k = m$ term.

- Since $X - Q$ is uncorrelated with all the W_m , the variance term decomposes as

$$\text{Var}(X - Q) + \text{Var} \left(\sum_{k=1}^n (\hat{a}_k - a_k) W_k \right) .$$

The first term we can't affect by choice of a_k , whereas we can zero the second term by setting $a_k = \hat{a}_k$ for all k , so we do that.

- Having chosen $a_k = \hat{a}_k$ for all k , we can zero the bias term by setting

$$b = \hat{b} = \mathbb{E}(X) - \sum_{k=1}^n \hat{a}_k \mathbb{E}(W_k) .$$

Summarize as follows.

Linear MMSE Estimation II: If W_k , $1 \leq k \leq n$, are uncorrelated random variables, the random variable \hat{X}_{opt} that minimizes (1) over all linear estimators $\hat{X} = \sum_{k=1}^n a_k W_k + b$

is

$$\begin{aligned}
\hat{X}_{\text{opt}} &= \sum_{k=1}^n \hat{a}_k W_k + \hat{b} \\
&= \sum_{k=1}^n \frac{\text{Cov}(X, W_k)}{\text{Var}(W_k)} W_k + \left(\mathbb{E}(X) - \sum_{k=1}^n \frac{\text{Cov}(X, W_k)}{\text{Var}(W_k)} \mathbb{E}(W_k) \right) \\
&= \sum_{k=1}^n \frac{\text{Cov}(X, W_k)}{\text{Var}(W_k)} (W_k - \mathbb{E}(W_k)) + \mathbb{E}(X) .
\end{aligned}$$

Again in this case the optimal estimator is unbiased and the optimal estimation error $\tilde{X} = X - \hat{X}_{\text{opt}}$ is uncorrelated with any linear combination of W_1, \dots, W_n .

Suppose finally that we want to estimate X linearly in terms of several random variables Y_1, \dots, Y_n that aren't necessarily uncorrelated. As usual, we want to minimize (1) using a linear estimator of the form

$$\hat{X} = \sum_{k=1}^n a_k Y_k + b .$$

One effective approach to the problem begins by generating uncorrelated random variables W_1, \dots, W_n that “span the same space” as Y_1, \dots, Y_n . You can do this using a procedure much like Gram-Schmidt orthogonalization of vectors or by diagonalizing the covariance matrix Σ of the Y_k , whose (i, j) -entry is $\text{Cov}(Y_i, Y_j)$ for $1 \leq i, j \leq n$. Having generated the W_k , you estimate X as above in terms of the uncorrelated random variables W_1, \dots, W_n and then rewrite that estimate in terms of Y_1, \dots, Y_n . The math works out so that the following result holds.

Linear MMSE Estimation III: If Y_1, \dots, Y_n are not necessarily uncorrelated random variables, the random variable \hat{X}_{opt} that minimizes (1) over all linear estimators $\hat{X} = \sum_{k=1}^n a_k Y_k + b$ is

$$\hat{X}_{\text{opt}} = (\text{Cov}(X, \mathbf{Y}))^T \Sigma^{-1} (\mathbf{Y} - \mathbb{E}(\mathbf{Y})) + \mathbb{E}(X) ,$$

where \mathbf{Y} is the column n -vector whose k th entry is Y_k , and $\text{Cov}(X, \mathbf{Y})$ is the column n -vector whose k th entry is $\text{Cov}(X, Y_k)$, and Σ is the $(n \times n)$ matrix whose (i, j) entry is $\text{Cov}(Y_i, Y_j)$, i.e. $\Sigma = \mathbb{E}(\mathbf{Y}\mathbf{Y}^T)$.