Modern Digital Communications: A Hands-On Approach

MMSE Estimation and LS Approximation

Dr. Nicolae Chiurtu

- Course material of Prof. Bixio Rimoldi -

Last revision: Nov. 27, 2023

Minimum Mean Squared Error Estimation

In this note, we study the concept of Minimum Mean Squared Error (MMSE) Estimation. We start with a general description of the problem, then focus on the special case of MMSE estimation for jointly Gaussian random vectors — which is what we need in order to estimate the OFDM channel.

Subsequently we study Least Squares (LS) approximations, which will also be used in the assignments.

Estimation

We are interested in estimating the realization x of the random variable $X \in \mathbb{C}^m$.

For this, we have access to the realization y of $Y \in \mathbb{C}^n$. (Here we use small letters to denote the realization of random variables). Unless X and Y are independent, knowing y should help estimating x.

An estimator of $X \in \mathbb{C}^m$ based on $Y \in \mathbb{C}^n$ can be an arbitrary function $\hat{X} : \mathbb{C}^n \to \mathbb{C}^m$. We then let $\hat{x} = \hat{X}(y)$ be the estimate of x based on the observed y. Clearly this estimate may or may not be a good one.

Thus we need a way to measure the performance of an estimator, and hopefully we can find an estimator which is optimal under that performance criterion.

MMSE Estimator

A popular choice for the fidelity criterion is the mean squared error

$$MSE(\hat{\boldsymbol{X}}) \coloneqq E \left[\left\| \boldsymbol{X} - \hat{\boldsymbol{X}}(\boldsymbol{Y}) \right\|^2 \right]$$

where the expectation is over the joint distribution of X and Y.

An estimator that minimizes the mean squared error is called a minimum $mean\ squared\ error\ (MMSE)\ estimator\ of\ {m X}\ given\ {m Y}.$

We denote by $\hat{m{X}}_{\mathrm{MMSE}}(m{y})$ such an estimator. In other words, for any estimator $\hat{m{X}}(m{y})$,

$$E\left[\left\|\boldsymbol{X} - \hat{\boldsymbol{X}}(\boldsymbol{Y})\right\|^{2}\right] \geq E\left[\left\|\boldsymbol{X} - \hat{\boldsymbol{X}}_{\mathrm{MMSE}}(\boldsymbol{Y})\right\|^{2}\right].$$

The MMSE estimator of \boldsymbol{X} given \boldsymbol{Y} turns out to be unique and equal to the $conditional\ expectation$ of \boldsymbol{X} given \boldsymbol{Y} , i.e.,

$$\hat{\boldsymbol{X}}_{\mathrm{MMSE}}(\boldsymbol{y}) = E\left[\boldsymbol{X}|\boldsymbol{Y}=\boldsymbol{y}\right].$$

Derivation of the MMSE Estimator

We are looking for an estimator $\hat{\boldsymbol{X}}\colon\mathbb{C}^n o\mathbb{C}^m$ such that

$$MSE(\hat{\boldsymbol{X}}) = E\left[\left\|\boldsymbol{X} - \hat{\boldsymbol{X}}(\boldsymbol{Y})\right\|^{2}\right],$$

is minimized.

We are done if we prove that $\hat{m{X}}_{ ext{MMSE}}(m{y})$ minimizes

$$E_y \left[\left\| \boldsymbol{X} - \hat{\boldsymbol{X}}(\boldsymbol{Y}) \right\|^2 \right] \stackrel{\triangle}{=} E \left[\left\| \boldsymbol{X} - \hat{\boldsymbol{X}}(\boldsymbol{Y}) \right\|^2 \middle| \boldsymbol{Y} = \boldsymbol{y} \right],$$

where E_y means that we are taking the expectation conditioning on $\boldsymbol{Y}=\boldsymbol{y}$.

We have

$$E_y \left[\left\| \boldsymbol{X} - \hat{\boldsymbol{X}}(\boldsymbol{Y}) \right\|^2 \right] = E_y \left[\left\| \boldsymbol{X} \right\|^2 \right] + \left\| \hat{\boldsymbol{X}}(\boldsymbol{y}) \right\|^2 - 2\Re \left\{ E_y \left[\boldsymbol{X}^{\dagger} \right] \hat{\boldsymbol{X}}(\boldsymbol{y}) \right\}.$$

The first term on the RHS of the equality does not depend on \hat{X} . If we substitute it for another term that does not depend on \hat{X} , then the minimizing \hat{X} remains the same.

We choose to minimize

$$\|E_y[\mathbf{X}]\|^2 + \|\hat{\mathbf{X}}(\mathbf{y})\|^2 - 2\Re \{E_y[\mathbf{X}^{\dagger}]\hat{\mathbf{X}}(\mathbf{y})\},$$

which is the quadratic form

$$||E_y[\boldsymbol{X}] - \hat{\boldsymbol{X}}(\boldsymbol{y})||^2.$$

The above is non-negative, and it achieves its minimum iff $\hat{\boldsymbol{X}}(\boldsymbol{y}) = E_y[\boldsymbol{X}]$.

MMSE Estimator for Jointly Gaussian Random Vectors

In general, finding an expression for E[X|Y=y] is not a trivial task.

However, when \boldsymbol{X} and \boldsymbol{Y} are jointly Gaussian random vectors, the MMSE estimator has a simple linear form.

In the rest of this lecture we restrict our discussion to jointly Gaussian random vectors \boldsymbol{X} and \boldsymbol{Y} .

First we assume that X and Y are zero-mean vectors. Later on, we will take non-zero means into account.

A Fact About Jointly Gaussian Random Vectors

If $X \in \mathbb{C}^m$ and $Y \in \mathbb{C}^n$ are zero-mean jointly Gaussian random vectors, we can always write

$$\boldsymbol{X} = A\boldsymbol{Y} + \boldsymbol{Z}$$

for some matrix $A \in \mathbb{C}^{m \times n}$ (that we will determine shortly) and a zero-mean Gaussian vector $\mathbf{Z} \sim \mathcal{N}_{\mathcal{C}}(\mathbf{0}, K_{\mathbf{Z}})$ that is independent of \mathbf{Y} .

Proof. For any matrix $A \in \mathbb{C}^{m \times n}$,

$$Z \coloneqq X - AY$$

is a (zero-mean) complex Gaussian vector. (By definition, a linear combination of jointly Gaussian vectors leads to a Gaussian vector.)

Moreover,

$$E\left[\mathbf{Z}\mathbf{Y}^{\dagger}\right] = E\left[\mathbf{X}\mathbf{Y}^{\dagger}\right] - AE\left[\mathbf{Y}\mathbf{Y}^{\dagger}\right] = K_{\mathbf{X}\mathbf{Y}} - AK_{\mathbf{Y}}$$

Hence, if we take

$$A = K_{\boldsymbol{X}\boldsymbol{Y}} K_{\boldsymbol{Y}}^{-1},$$

we will have

$$E\left[\mathbf{Z}\mathbf{Y}^{\dagger}\right]=0,$$

i.e., with this specific choice of A , \boldsymbol{Z} is independent of \boldsymbol{Y} .

Note that with this choice of A,

$$K_{\mathbf{Z}} = E \left[(\mathbf{X} - A\mathbf{Y})(\mathbf{X} - A\mathbf{Y})^{\dagger} \right]$$

$$= K_{\mathbf{X}} - AK_{\mathbf{Y}\mathbf{X}} - (K_{\mathbf{X}\mathbf{Y}} - AK_{\mathbf{Y}})A^{\dagger}$$

$$= K_{\mathbf{X}} - AK_{\mathbf{Y}\mathbf{X}}$$

$$= K_{\mathbf{X}} - K_{\mathbf{X}\mathbf{Y}}K_{\mathbf{Y}}^{-1}K_{\mathbf{X}\mathbf{Y}}^{\dagger},$$

where in the third line we used the fact that $AK_{Y} = K_{XY}$.

MMSE Estimator for Jointly Gaussian Vectors

Let $oldsymbol{X}$ and $oldsymbol{Y}$ be jointly Gaussian. Then

$$\hat{\boldsymbol{X}}_{\mathrm{MMSE}}(\boldsymbol{y}) = K_{\boldsymbol{X}\boldsymbol{Y}}K_{\boldsymbol{Y}}^{-1}\boldsymbol{y}.$$

Moreover, the minimum mean squared error equals

$$E\left[\left\|\boldsymbol{X} - \hat{\boldsymbol{X}}_{\text{MMSE}}(\boldsymbol{Y})\right\|^{2}\right] = \operatorname{trace}\left(K_{\boldsymbol{X}} - K_{\boldsymbol{X}\boldsymbol{Y}}K_{\boldsymbol{Y}}^{-1}K_{\boldsymbol{X}\boldsymbol{Y}}^{\dagger}\right).$$

Proof. Write

$$\boldsymbol{X} = A\boldsymbol{Y} + \boldsymbol{Z}$$
 with $A = K_{\boldsymbol{X}\boldsymbol{Y}}K_{\boldsymbol{Y}}^{-1}$,

which makes $oldsymbol{Z}$ independent of $oldsymbol{Y}$.

Now

$$\hat{\boldsymbol{X}}_{\mathrm{MMSE}}(\boldsymbol{y}) = E\left[\boldsymbol{X}|\boldsymbol{Y} = \boldsymbol{y}\right]$$

= $E\left[A\boldsymbol{Y} + \boldsymbol{Z}|\boldsymbol{Y} = \boldsymbol{y}\right]$
= $E\left[A\boldsymbol{y} + \boldsymbol{Z}\right] = A\boldsymbol{y}.$

The minimum mean squared error is

$$E\left[\left\|\boldsymbol{X} - \hat{\boldsymbol{X}}_{\text{MMSE}}(\boldsymbol{Y})\right\|^{2}\right] = E\left[\left\|(A\boldsymbol{Y} + \boldsymbol{Z}) - A\boldsymbol{Y}\right\|^{2}\right]$$
$$= E\left[\left\|\boldsymbol{Z}\right\|^{2}\right] = \operatorname{trace}(K_{\boldsymbol{Z}}) = \operatorname{trace}(K_{\boldsymbol{X}} - K_{\boldsymbol{X}\boldsymbol{Y}}K_{\boldsymbol{Y}}^{-1}K_{\boldsymbol{X}\boldsymbol{Y}}^{\dagger}).$$

Channel Coefficients Estimation

We have seen that by using OFDM we 'forge' parallel channels described in matrix form as

$$\boldsymbol{Y}^{(m)} = D\boldsymbol{A}^{(m)} + \boldsymbol{Z}^{(m)}$$

where D is the diagonal matrix of channel coefficients. How to estimate the matrix D?

For certain values of m, we substitute $A^{(m)}$ with an N-tuple S known to the receiver. Then, dropping the superscript (m) for notational convenience,

$$Y = DS + Z$$

where D is diagonal. The same result is obtained by

$$Y = SD + Z$$

where S is the diagonal matrix that has S as its diagonal elements and D is the N-tuple consisting of the diagonal elements of D.

Assume that $oldsymbol{D}$ is zero-mean, Gaussian, and independent of $oldsymbol{Z}$.

Then $oldsymbol{Y}$ and $oldsymbol{D}$ are jointly Gaussian. To see this, write

$$\begin{pmatrix} \boldsymbol{Y} \\ \boldsymbol{D} \end{pmatrix} = \begin{pmatrix} S & I \\ I & 0 \end{pmatrix} \begin{pmatrix} \boldsymbol{D} \\ \boldsymbol{Z} \end{pmatrix}.$$

This shows that the vector on the left hand side is a linear transformation of a Gaussian random vector.

Hence, the minimum mean squared error (MMSE) estimate of $m{D}$ based on the observable $m{Y}$ is

$$\hat{\boldsymbol{D}} = K_{\boldsymbol{D}\boldsymbol{Y}} K_{\boldsymbol{Y}}^{-1} \boldsymbol{Y},$$

where

$$K_{DY} := E[DY^{\dagger}] = K_{D}S^{\dagger}$$
 and $K_{Y} := E[YY^{\dagger}] = SK_{D}S^{\dagger} + K_{Z}$

are covariance matrices and $oldsymbol{D}$ and $oldsymbol{Y}$ are zero-mean.

Both K_{Y} and K_{DY} depend on K_{D} . The question is how to find K_{D} .

K_D from the Channel Model

To determine $K_{\mathbf{D}}$, we make a channel model.

A reasonable assumption for a wireless channel is the multipath channel, modeled as

$$h(t) = \sum_{l=0}^{M-1} \alpha_l \delta(t - \tau_l),$$

where α_l and τ_l are the *l*-th path strength and delay, respectively.

Since α_l and α_k , $l \neq k$, describe the reflection on different obstacles/surfaces, it is reasonable to assume that they are uncorrelated.

For this channel,

$$h_{\mathcal{F}}(f) = \int h(t)e^{-j2\pi ft}dt = \sum_{l} \alpha_{l} \int \delta(t - \tau_{l})e^{-j2\pi ft}dt = \sum_{l} \alpha_{l}e^{-j2\pi f\tau_{l}}.$$

Assuming that the DFT-length N is even (usually a power of 2), $\mathbf{D} = (\lambda_0, \dots, \lambda_{N-1})$ with

$$\lambda_i = \begin{cases} h_{\mathcal{F}} \left(\frac{i}{NT_s} \right), & i = 0, \dots, \frac{N}{2} - 1 \\ h_{\mathcal{F}} \left(\frac{i-N}{NT_s} \right), & i = \frac{N}{2}, \dots, N - 1 \end{cases},$$

where

$$h_{\mathcal{F}}\left(\frac{i}{NT_s}\right) = \sum_{l} \alpha_l e^{-j2\pi \frac{i}{NT_s}\tau_l},$$

and T_s is both the symbol interval and the sampling interval.

If we define

$$[i] = \begin{cases} i, & i = 0, \dots, \frac{N}{2} - 1 \\ i - N, & i = \frac{N}{2}, \dots, N - 1 \end{cases}$$

then we can write

$$\lambda_i = h_{\mathcal{F}} \left(\frac{[i]}{NT_s} \right).$$

The (i, k) entry of the covariance matrix $K_{\mathbf{D}}$ is

$$(K_{\mathbf{D}})_{i,k} = E\left[\sum_{l} \alpha_{l} e^{-j2\pi \frac{[i]}{NT_{S}} \tau_{l}} \sum_{v} \alpha_{v}^{*} e^{j2\pi \frac{[k]}{NT_{S}} \tau_{v}}\right]$$

$$= \sum_{l} \sum_{v} e^{-j2\pi \frac{[i]}{NT_{S}} \tau_{l}} E\left[\alpha_{l} \alpha_{v}^{*}\right] e^{j2\pi \frac{[k]}{NT_{S}} \tau_{v}}$$

$$(1)$$

The above expression seems complicated but it is actually not. In fact, K_D is the product of three matrices. To see this, notice that if A,B,C are matrices such that the product ABC is well defined, then

$$(ABC)_{i,k} = \sum_{l} \sum_{v} A_{i,l} B_{l,v} C_{v,k},$$

which is exactly the right-hand side of (1) with

$$A_{i,l} = e^{-j2\pi \frac{[i]}{NT_s}\tau_l}, \qquad B_{l,v} = E\left[\alpha_l \alpha_v^*\right], \qquad C = A^{\dagger}.$$

From $K_{\mathbf{D}}$, we determine $K_{\mathbf{Y}}$ and $K_{\mathbf{DY}}$ as described earlier.

Least Squares (LS) Approximation

An MMSE estimator requires the knowledge of the statistics.

Suppose that all we know is that the observable $m{y} \in \mathbb{C}^n$ is obtained from $m{\lambda} \in \mathbb{C}^n$ according to

$$y = S\lambda + z$$

where S is a diagonal matrix with non-vanishing diagonal elements, and $z \in \mathbb{C}^n$ is a noise vector. None of the statistics are known.

We would like to find the $\hat{\lambda} \in \mathbb{C}^n$ for which $\|y - S\lambda\|^2$ is minimized over all $\lambda \in \mathbb{C}^n$.

By writing

$$\|y - S\lambda\|^2 = \sum_{i=0}^{n-1} |y_i - S_i\lambda_i|^2,$$

it is clear that $\hat{\lambda}_i$ needs to be chosen to minimize $|y_i - S_i \lambda_i|^2$.

Clearly the minimum is obtained when

$$\hat{\lambda}_i = \frac{y_i}{S_i},$$

in which case $\|\boldsymbol{y} - S\boldsymbol{\lambda}\|^2 = 0$.

This is a special case that we have worked out since we will need it in the assignment.

More generally, suppose that

$$y = S\lambda + z$$

where y and z are in \mathbb{C}^n , $\lambda \in \mathbb{C}^m$, and $S \in \mathbb{C}^{n \times m}$ is a general full-rank matrix (not necessarily diagonal).

As before, we are seeking the $\hat{\lambda}$ that minimizes

$$\|\boldsymbol{y} - S\boldsymbol{\lambda}\|^2$$

over all $\lambda \in \mathbb{C}^m$.

If $n \leq m$, the system $\mathbf{y} = S\boldsymbol{\lambda}$ can always be solved. Hence we can always find a $\boldsymbol{\lambda}$ for which $\|\mathbf{y} - S\boldsymbol{\lambda}\|^2 = 0$.

So we focus on the case where n > m (the system is overdetermined).

We may reformulate the problem as follows. The observable y is an element of the inner-product space $\mathcal{U} = \mathbb{C}^n$ and let \mathcal{V} be the subspace spanned by the columns of S. We are seeking the vector $\hat{y} \in \mathcal{V}$ that minimizes

$$\|\boldsymbol{y} - \hat{\boldsymbol{y}}\|^2$$
.

The projection theorem (see PDC) tells us that \hat{y} is the projection of y into \mathcal{V} . It has the property that the error vector $y - \hat{y}$ is orthogonal to every element of \mathcal{V} . In particular, it is orthogonal to the columns of S. Hence

$$\langle \boldsymbol{y} - S\hat{\boldsymbol{\lambda}}, \boldsymbol{S}_i \rangle = 0, \quad i = 1, \dots, m$$

where S_i is the *i*-th column of S.

Hence

$$\langle \boldsymbol{y}, \boldsymbol{S}_i \rangle = \langle S \hat{\boldsymbol{\lambda}}, \boldsymbol{S}_i \rangle, \quad i = 1, \dots, m.$$
 (2)

Equivalently,

$$\boldsymbol{S}_{i}^{\dagger}\boldsymbol{y}=\boldsymbol{S}_{i}^{\dagger}S\hat{\boldsymbol{\lambda}}\quad i=1,\ldots,m.$$

In matrix form, we obtain the so-called *normal equations*:

$$S^{\dagger} \boldsymbol{y} = S^{\dagger} S \hat{\boldsymbol{\lambda}}.$$

Since S has full rank, $S^{\dagger}S$ is nonsingular. We prove this by arguing that if $S^{\dagger}S$ is singular, then S does not have full rank. Indeed, if $S^{\dagger}S$ is singular, there exists a nonzero vector \boldsymbol{u} such that $S^{\dagger}S\boldsymbol{u}=0$, which implies that $\boldsymbol{u}^{\dagger}S^{\dagger}S\boldsymbol{u}=\|S\boldsymbol{u}\|^2=0$, so that $S\boldsymbol{u}=0$. This means that the columns of S are linearly dependent, i.e., S is not full rank — a contradiction.

Solving for $\hat{\lambda}$ yields the least-squares approximation of λ :

$$\hat{\boldsymbol{\lambda}} = \left(S^{\dagger} S \right)^{-1} S^{\dagger} \boldsymbol{y}.$$

To reconstruct the above formula from memory, it suffices to check the dimensions. The matrix that estimates λ from y must be of dimension $m \times n$. If we analyze $(S^\dagger S)^{-1} S^\dagger$ from right to left, we see that we don't have other choices (unless we make the expression more complicated). In particular, S^\dagger has the right dimension to multiply y and notice that $(S^\dagger S)^{-1}$ is $m \times m$, hence it is a valid matrix to multiply the $m \times n$ matrix S^\dagger , whereas $(SS^\dagger)^{-1}$ would not fit as it is $n \times n$.

To summarize, recall the fundamental difference between the MMSE estimate and the LS approximation. On the one hand, the MMSE setup assumes two random vectors, the unobserved \boldsymbol{X} and the observed \boldsymbol{Y} , and the objective is to estimate \boldsymbol{X} by means of \boldsymbol{Y} , where the measure of performance is $E[\|\boldsymbol{X} - \hat{\boldsymbol{X}}(\boldsymbol{Y})\|^2]$.

On the other hand, the LS approximation is about adjusting the parameters of a model function to best fit the observed data. (See the example that follows.) It is also called data fitting or regression analysis. (See e.g. Least Squares in Wikipedia).

Example: Suppose that the model function is $y(x) = \alpha_1 + \alpha_2 x + \alpha_3 x^2$, with unknown parameters α_1 , α_2 , α_3 . We seek to find the parameters so that the model fits the noisy observations

$$y_i = y(x_i) + z_i, \quad i = 1, 2, \dots, n.$$

Letting $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$, $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$, $\mathbf{\alpha} = (\alpha_1, \alpha_2, \alpha_3)^T$, and $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$, the observations are of the form

$$y = S\alpha + z,$$

where

$$S = \begin{pmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix}.$$

It is worth pointing out that even if we are estimating a nonlinear function of x, we are dealing with a data observation model where x has fixed values and the variable is the parameter vector α . As a function of α , the observation model is linear.

The $\pmb{\alpha}$ that minimizes $\|\pmb{y}-\hat{\pmb{y}}(\pmb{x})\|^2$ is the LS approximation $\hat{\pmb{\alpha}}=(S^\dagger S)^{-1}S^\dagger \pmb{y}.$

For a MATLAB example, check out the file example.m