

Stochastic Processes

Session 16

Minimum Mean Square Error Estimation

- 1) Recap Estimation, Errors, Bias & MSE
- 2) MMSE
- 3) LMMSE
- 4) Vector LMMSE

Estimation Problem

Given $X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$ estimate θ .

Exercise: State a few estimation problems relevant to your field of study

Estimator

Def: An estimator is a function $\hat{\theta} = g(X)$ of data X used to guess the value of some unknown entity θ .

Recall the following defs.:

- Estimation error: $\theta - \hat{\theta}$
- Bias (mean error): $E[\theta - \hat{\theta}]$
- Mean Squared Error (MSE):
$$E[(\theta - \hat{\theta})^2]$$

Exercise: Suppose we ignore the data X and just guess on the mean: $\hat{\theta} = E[\theta] = \mu_{\theta}$. What is the bias and the MSE of this estimator?

Bayesian estimation

- We assume θ to be a random variable with a priori pdf $p(\theta)$
 - θ is jointly distributed with the data X according to $p(\theta, X)$.
-
- A priori knowledge may be obtained from past experiments or information on the context.
 - In general the a priori information may be full (i.e. $p(\theta)$) or partial (e.g. moments of $p(\theta)$).

Exercise: Can you phrase your estimation problem from before as a Bayesian Estimation problem? What a priori information can be used in this case?

Minimum MSE estimator (MMSE)

Def.: The MMSE of θ given the data X is the function $\hat{\theta} = g(X)$ such that the MSE $E[(\theta - \hat{\theta})^2]$ is minimum.

Notice that $E[(\theta - \hat{\theta})^2] = E[E[(\theta - \hat{\theta})^2 | X]]$.

Therefore, $\hat{\theta}$ should be chosen to minimize the conditional MSE $E[(\theta - \hat{\theta})^2 | X]$ for each X .

We find $\hat{\theta}$ by "differentiating and equating to zero". The chain-rule of differentiation simplifies this task.

Derivation of the MMSE

It suffices to minimize the conditional MSE wrt $\hat{\theta}$:

$$\frac{\partial E[(\theta - \hat{\theta})^2 | X]}{\partial \hat{\theta}} = E \left[\frac{\partial (\theta - \hat{\theta})^2}{\partial \hat{\theta}} | X \right]$$

$$= E \left[2(\theta - \hat{\theta}) \cdot \frac{\partial (\theta - \hat{\theta})}{\partial \hat{\theta}} | X \right]$$

$$= -2 E[\theta - \hat{\theta} | X]$$

$$= -2 E[\theta | X] + 2 \hat{\theta}$$

(Since $\hat{\theta} = g(X)$ is "deterministic" cond. on X)

Equating to zero, we obtain:

$$\boxed{\hat{\theta} = E[\theta | X]}$$

(Double diff. shows that this critical point is indeed minimum)

Properties of the MMSE:

- The MMSE is unbiased:
$$E[\theta - \hat{\theta}] = E[\theta - E[\theta|X]] = E[\theta] - E[\theta] = 0$$
- The MMSE has the lowest MSE among all estimators.
- The MMSE fulfills the "orthogonality principle", i.e. the error is uncorr. to any function of the data:

$$E[(\theta - \hat{\theta}) \cdot h(X)] = 0$$

Proof:
$$E[(\theta - \hat{\theta}) \cdot h(X)] = E[E[(\theta - \hat{\theta})h(X)|X]]$$
$$= E[\underbrace{E[\theta - \hat{\theta}|X]}_{=0} \cdot h(X)] = 0.$$

(see deriv of MMSE)

Remarks to the MMSE

- Although simple to derive, the MMSE is impractical unless $E[\theta|X]$ can be computed
- Analytic solutions are available only in (very) special cases.
- Alternatively $E[\theta|X]$ may be computed numerically (e.g. via Monte Carlo simulation), or using other approximation techniques.
- Instead of the MMSE, a linear MMSE (LMMSE) may be applied.

Linear MMSE (LMMSE)

Def: The LMMSE is the estimator of the form

$$\hat{\theta} = h_0 + \sum_{n=1}^N h_n X_n = h_0 + \mathbf{h}^T \mathbf{X}$$

where h_0, h_1, \dots, h_N are chosen to minimize $E[(\theta - \hat{\theta})^2]$.

Remark: The estimator is actually not linear, but "affine" in cases where $h_0 \neq 0$. However, the name is used widely in the literature so we stick to the convention.

• The coefficients h_0, \dots, h_N can be obtained directly by minimizing the MSE.

Coeffs of the LMMSE

As with the MMSE, we make use of the chain-rule to minimize $E[(\theta - \hat{\theta})^2 | X]$:

$$\begin{aligned} \bullet \quad \frac{\partial E[(\theta - \hat{\theta})^2 | X]}{\partial h_0} &= \underbrace{\frac{\partial \hat{\theta}}{\partial h_0}}_{=1} \cdot \underbrace{\frac{\partial E[(\theta - \hat{\theta})^2 | X]}{\partial \hat{\theta}}}_{-2 E[\theta - \hat{\theta} | X]} \\ &= -2 E[\theta | X] + 2(h_0 + h^T X) \end{aligned}$$

Equating to zero and taking the expectation w.r.t. X gives

$$0 = -E[\theta] + h_0 + h^T E[X] \Rightarrow \underline{\underline{h_0 = E[\theta] - h^T E[X]}}$$

$$\begin{aligned} \bullet \quad \frac{\partial E[(\theta - \hat{\theta})^2 | X]}{\partial h^T} &= \underbrace{\frac{\partial \hat{\theta}}{\partial h^T}}_{=X} \cdot \underbrace{\frac{\partial E[(\theta - \hat{\theta})^2 | X]}{\partial \hat{\theta}}}_{-2 E[\theta - \hat{\theta} | X]} \\ &= -2 E[(\theta - \hat{\theta}) \cdot X | X] \end{aligned}$$

Equating to zero and taking the expectation w.r.t. X gives

$$E[(\theta - (h_0 + h^T X)) \cdot X] = 0$$

→

Inserting for $h_0 \Rightarrow$

$$E[\underbrace{(\theta - E[\theta] - h^T (X - E[X])) \cdot X}_{\text{zero mean}}] = 0$$

we notice that this is zero mean, and thus we see that

$$E[(\theta - E[\theta] - h^T (X - E[X])) \cdot (X - E[X])] = 0$$

$$\Downarrow$$
$$C_{\theta X} - h^T C_{XX} = 0$$

\Downarrow

$$h^T = C_{\theta X} C_{XX}^{-1}$$

\Downarrow

$$\underline{\underline{h = C_{XX}^{-1} C_{X\theta}}}$$

$$(C_{\theta X}^T = C_{X\theta})$$

The LMMSE now reads

$$\hat{\theta} = E[\theta] - C_{\theta X} C_{XX}^{-1} E[X] + C_{\theta X} C_{XX}^{-1} X.$$

Alternative form:

$$\hat{\theta} = E[\theta] + C_{\theta X} C_{XX}^{-1} (X - E[X])$$

Remarks

- The LMMSE requires a priori knowledge of $E[\theta]$, $E[X]$, $C_{\theta X}$ and C_{XX} .
- Simple to compute provided C_{XX} is invertible

Properties of the LMMSE

- The LMMSE is unbiased
(see derivation of h_0)
- The LMMSE minimizes the MSE among linear (or rather affine) estimators.
 - Nonlinear estimators may perform better!
- The LMMSE fulfills the "orthogonality principle":

$$E[(\theta - \hat{\theta}) f(x)] = 0$$

where $f(x)$ is any affine function of X , i.e. $f(x) = f_0 + \underline{f}^T \underline{x}$.

Proof:
$$E[(\theta - \hat{\theta}) f(x)] = \underbrace{E[(\theta - \hat{\theta}) f_0]}_{=0, \hat{\theta} \text{ unbiased.}} + E[(\theta - \hat{\theta}) \underline{f}^T \underline{x}]$$
$$= \underline{f}^T \underbrace{E[(\theta - \hat{\theta}) \underline{x}]}_{=0 \text{ see deriv of } \underline{h}.}$$

Residual MSE of LMMSE

The residual MSE can be computed by application of the orthogonality principle (o.p.):

$$\begin{aligned} E[(\theta - \hat{\theta})^2] &= E[(\theta - \hat{\theta})(\theta - \hat{\theta})] \\ &= E[(\theta - \hat{\theta})\theta] \quad (\text{o.p.}) \\ &= E[(\theta - \hat{\theta})(\theta - E[\theta])] \quad (\hat{\theta} \text{ unbiased}) \\ &= E\left[\left((\theta - E[\theta]) - C_{\theta x} C_{xx}^{-1} (x - E[x])\right)(\theta - E[\theta])\right] \\ &\quad (\text{alt. form f. } \hat{\theta}) \end{aligned}$$

$$\underline{\underline{= \text{Var}(\theta) - C_{\theta x} C_{xx}^{-1} C_{x\theta}.}}$$

- So the MSE of the LMMSE is known upfront!

LMMSE for multiple variables ("Vector LMMSE")

The unknown parameter $\underline{\theta}$ is now a vector with K entries:

$$\underline{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_K \end{bmatrix}$$

- We seek the linear (or affine) estimator $\hat{\underline{\theta}}$ which minimizes the total MSE:

$$MSE_{tot} = MSE_1 + MSE_2 + \dots + MSE_K$$

$$\text{where } MSE_k = E[(\theta_k - \hat{\theta}_k)^2].$$

- MSE_{tot} is minimized if and only if MSE_1, \dots, MSE_K are all minimum.

\Rightarrow We simply have to use parallel scalar LMMSEs !

Expression for "Vector LMMSE"

We stack scalar LMMSEs in a vector:

$$\hat{\underline{\theta}} = \begin{bmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_K \end{bmatrix} = \begin{bmatrix} E[\theta_1] + C_{\theta_1 X} C_{XX}^{-1} (X - E[X]) \\ \vdots \\ E[\theta_K] + C_{\theta_K X} C_{XX}^{-1} (X - E[X]) \end{bmatrix}$$

$$= E[\theta] + C_{\theta X} C_{XX}^{-1} (X - E[X])$$

The total MSE is then:

$$MSE_{tot} = \sum_{k=1}^K MSE_k$$

$$= \sum_{k=1}^K \left(\text{Var}(\theta_k) - C_{\theta_k X} C_{XX}^{-1} C_{X \theta_k} \right)$$

$$= \text{trace} (C_{\theta\theta} - C_{\theta X} C_{XX}^{-1} C_{X\theta}).$$