→ 3rd Lecture : Tuesday Jan 16th 2018

- **M̲S̲E̲ :**

To study estimation error we started by studying $P(|\hat{\theta}_n - \theta| > \delta)$, deviation above a given threshold $\delta$, by bounding this probability. One may take a different approach by studying average Euclidean distance, i.e. $E[|\hat{\theta}_n - \theta|^2]$, which is denoted by $MSE(\hat{\theta}_n)$.

We note that if $\theta = E(\hat{\theta}_n)$, i.e. $\hat{\theta}_n$ is an unbiased estimate of $\theta$, then $MSE(\hat{\theta}_n) = E[|\hat{\theta}_n - \theta|^2] = E[(\hat{\theta}_n - \underset{\hat{\theta}_n}{\mu})^2] = Var(\hat{\theta})$.

Now recall that $Var(X) = 0 \Rightarrow P(X = constant) = 1$ which essentially means r.X. is a constant. The same comment applies to $MSE(\hat{\theta}_n)$. We want to find the closest estimator $\hat{\theta}_n$ to $\theta$ which means that we want to minimize $E[(\hat{\theta}_n - \theta)^2]$ over all possible estimators, ideally attain

the above comment tell us that in real applications we cannot expect to find an estimator whose MSE is equal to zero. Let's try to understand the MSE a bit more

$$MSE(\hat{\theta}_n) = E[(\hat{\theta}_n - \theta)^2]$$
$$= E[\{(\hat{\theta}_n - E(\hat{\theta}_n)) + (E(\hat{\theta}_n) - \theta)\}^2]$$
$$= E[\{\hat{\theta}_n - E(\hat{\theta}_n)\}^2 + \{E(\hat{\theta}_n) - \theta\}^2 + 2\{E(\hat{\theta}_n) - \theta\}\{\hat{\theta}_n - E(\hat{\theta}_n)\}]$$
$$= E[\{\hat{\theta}_n - E(\hat{\theta}_n)\}^2] + E[\{\underbrace{E(\hat{\theta}_n) - \theta}_{NoT\ a\ r.v.}\}^2]$$
$$+ 2E[\{\underbrace{E(\hat{\theta}_n) - \theta}_{NoT\ a\ r.v.}\}\{\hat{\theta}_n - E(\hat{\theta}_n)\}]$$

$$= \text{Var}\left(\hat{\theta}_n\right) + \underbrace{\left[E(\hat{\theta}_n) - \theta\right]^2}_{\text{Bias}(\hat{\theta}_n)} + 2\,\text{Bias}(\hat{\theta}_n)\,\underbrace{E\left[(\hat{\theta}_n - E(\hat{\theta}_n))\right]}_{\substack{0 \\ E(\hat{\theta}_n) - E(\hat{\theta}_n) =}}$$

$$= \text{Var}(\hat{\theta}_n) + \text{Bias}^2(\hat{\theta}_n)$$

Roughly speaking, bias measures how far off the target we hit on the average while variance measures how much fluctuation our estimator may show from one sample to another.

- Unbiased Estimators:

  In almost all real applications, the class of possible estimators for an estimand is huge and the best estimator, i.e. the one that minimizes MSE no matter what the value of the estimand is, almost never exists. Thus we try to reduce the class of potential estimators by imposing a plausible restriction, for example Bias$(\hat{\theta}_n) = 0$

- Def. An estimator $\hat{\theta}_n$ of an estimand $\theta$ is said to be unbiased if $E(\hat{\theta}_n) = \theta$, for all possible values of.

  Example: $X_i \overset{iid}{\sim} N(\mu, \sigma^2)$
  $\qquad\qquad i = 1, 2, -, n$
  suppose both $\mu$ and $\sigma^2$ are unknown. Consider
  $\bar{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i$.

  $$E(\bar{X}_n) = E\left(\frac{1}{n}\sum_{i=1}^{n} X_i\right) = \frac{1}{n}\sum_{i=1}^{n} E(\overset{\mu}{\overbrace{X_i}}) = \frac{1}{\not{n}} \cdot \not{n}\,\mu = \mu.$$

  Thus $\bar{X}_n$ is an unbiased estimator of $\mu$. As for the MSE$(\bar{X}_n)$, we need to find Var$(\bar{X}_n)$.

$$\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n}\sum_{i=1}^{n}X_i\right) = \frac{1}{n^2}\text{Var}\left(\sum_{i=1}^{n}X_i\right)$$

Thm 5.12(b)
page 271
$$= \frac{1}{n^2}\left\{\sum_{i=1}^{n}\text{Var}(X_i) + 2\sum_{1\le i<j\le n}\underbrace{\text{Cov}(X_i, X_j)}_{0}\right\}$$

$\overset{n}{\underset{i=1}{\amalg}} X_i$
$$= \frac{1}{n^2}\sum_{i=1}^{n}\text{Var}(X_i)$$

Identically
distributed
$$= \frac{1}{n^2}\sum_{i=1}^{n}\sigma^2 = \frac{1}{n\!\!\!/}\times n\!\!\!/\times\sigma^2 = \frac{\sigma^2}{n}$$

$$\text{MSE}(\bar{X}_n) = \text{Var}(\bar{X}_n) + \underbrace{\text{Bias}^2(\bar{X}_n)}_{0} = \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$$

An inspection of the above calculation shows that for unbiased ness we only require a common mean $\mu$ while for calculating the variance we would only require a common variance $\sigma^2$ and orthogonality, i.e.
$$\text{Cov}(X_i, X_j) = 0 \quad \text{if } i \ne j.$$

— suppose $X_1, \dots, X_n$ have the same mean value $\mu$. Then
$$E(\bar{X}_n) = E\left(\frac{1}{n}\sum_{i=1}^{n}X_i\right) = \frac{1}{n}\sum_{i=1}^{n}E(X_i) = \frac{1}{n\!\!\!/}\times n\!\!\!/\times\mu = \mu.$$

— suppose further that $X_1, \dots, X_n$ have the same varian $\sigma^2$ and $\text{Cov}(X_i, X_j) = 0$, $i \ne j$. Then
$$\text{Var}(\bar{X}_n) = \text{Var}\left(\frac{1}{n}\sum_{i=1}^{n}X_i\right) = \frac{1}{n^2}\text{Var}\left(\sum_{i=1}^{n}X_i\right)$$

Thm 5.12(b)
page 271
$$= \frac{1}{n^2}\left\{\sum_{i=1}^{n}\text{Var}(X_i) + 2\sum_{1\le i<j\le n}\text{Cov}(X_i, X_j)\right\}$$

orthogonality, i.e.
$\text{Cov}(X_i, X_j) = 0$
$i \ne j$
$$= \frac{1}{n^2}\sum_{i=1}^{n}\text{Var}(X_i)$$

having the same variance $= \frac{1}{n^2}\sum_{i=1}^{n}\sigma^2 = \frac{1}{n\!\!\!/}\times n\!\!\!/\times\sigma^2 = \frac{\sigma^2}{n}$

Thus $MSE(\bar{X}_n) = Var(\bar{X}_n) = \frac{\sigma^2}{n}$

if $X_1, \ldots, X_n$ have the same mean value and variance and they are orthogonal.


— Remark (Stein's Paradox)

We will learn later that if $X_i \overset{iid}{\sim} N(\mu, \sigma^2)$, $i=1,\ldots,n$ then $\bar{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i$ has many optimal properties. A paradox due to Charls Stein, however, shows that such nice optimal properties are not preserved in higher dimensions. In fact if

$$X_i \overset{iid}{\sim} N(\mu_x, 1), \quad Y_i \overset{iid}{\sim} N(\mu_Y, 1) \quad \text{and}$$
$$Z_i \overset{iid}{\sim} N(\mu_2, 1),$$

then we can find biased estimators of $\begin{pmatrix}\mu_x\\\mu_Y\\\mu_2\end{pmatrix}$ which are closer to $\begin{pmatrix}\mu_x\\\mu_Y\\\mu_2\end{pmatrix}$ than $\begin{pmatrix}\bar{X}_n\\\bar{Y}_n\\\bar{Z}_n\end{pmatrix}$ for any $\begin{pmatrix}\mu_x\\\mu_Y\\\mu_2\end{pmatrix}$

We may then say that $\begin{pmatrix}\bar{X}_n\\\bar{Y}_n\\\bar{Z}_n\end{pmatrix}$ is an

<u>inadmissible</u> estimator of $\begin{pmatrix}\mu_x\\\mu_Y\\\mu_2\end{pmatrix}$.

○ Admissibility:

An estimator $\hat{\theta}$ is called admissible if there is no estimator $\tilde{\theta}$ such that

$$MSE(\tilde{\theta}) \leq MSE(\hat{\theta}) \quad \text{for all possible values}$$
$$\text{of } \theta$$

and the inequality is strict for some values of $\theta$.

What this example tells us is that by allowing a bit of bias we may be able to reduce variance considerably and hence find an estimator which closer to the target than the most natural unbiased estimator. Note that this phenomenon happens only when the dimension is at least 3. (12)

- We now want to restrict the class of estimators even further. Suppose $X_1, -, X_n$ have the same mean $\mu$ and variance $\sigma^2$ and they are orthogonal, i.e $Cov(X_i, X_j) = 0$, $i \neq j$. Consider $\tilde{X}_{n,\underset{\sim}{c}} = \sum_{i=1}^{n} c_i X_i$ and

$$\mathcal{C} = \left\{ \tilde{X}_{n,\underset{\sim}{c}} : \underset{\sim}{c} = (c_1, -, c_n) \in \mathbb{R}^n, \sum_{i=1}^{n} c_i = 1 \right\}.$$

Note that

$$E(\tilde{X}_{n,\underset{\sim}{c}}) = E\left( \sum_{i=1}^{n} c_i X_i \right) = \sum_{i=1}^{n} c_i E(X_i)$$

$$= \sum_{i=1}^{n} c_i \mu = \mu \underbrace{\sum_{i=1}^{n} c_i}_{1} = 1 \times \mu = \mu$$

Thus $\tilde{X}_{n,\underset{\sim}{c}}$ is an unbiased estimator of $\mu$ for any $\underset{\sim}{c} \in \mathbb{R}^n$ as long as $\sum_{i=1}^{n} c_i = 1$. Then $\mathcal{C}$ is the class of all unbiased linear estimators of $\mu$. We want to find the best estimator within $\mathcal{C}$, i.e.

$$\underset{\underset{\sim}{c} \in \mathbb{R}^n}{\text{Min}} \text{ MSE}(\tilde{X}_{n,\underset{\sim}{c}}) \qquad (\dagger)$$

$$\text{s.t. } \sum_{i=1}^{n} c_i = 1$$

First we note that $\text{MSE}(\tilde{X}_{n,\underset{\sim}{c}}) = \text{Var}(\tilde{X}_{n,\underset{\sim}{c}})$ since $\tilde{X}_{n,\underset{\sim}{c}}$ is an unbiased estimator of $\mu$ when $\sum_{i=1}^{n} c_i = 1$. On the other hand

$$\text{Var}(\tilde{X}_{n,\underset{\sim}{c}}) = \text{Var}\left( \sum_{i=1}^{n} c_i X_i \right)$$

Thm 5.12
page 271

$$= \sum_{i=1}^{n} c_i^2 \text{Var}(X_i) + 2 \sum\sum_{1 \leq i < j \leq n} Cov(c_i X_i, c_j X_j)$$

$$= \sum_{i=1}^{n} c_i^2 \sigma^2 + 2 \sum\sum_{1 \leq i < j \leq n} c_i c_j \underbrace{Cov(X_i, X_j)}_{0}$$

$$= \sigma^2 \sum_{i=1}^{n} c_i^2$$

Thus (†) is equivalent to

$$\underset{\underset{\sim}{c} \in \mathbb{R}^n}{\text{Min}} \quad \sigma^2 \sum_{i=1}^{n} c_i^2 \qquad (\mp)$$

$$\text{s. t.} \quad \sum_{i=1}^{n} c_i = 1$$

using Lagrange Theorem (∓) is equivalent to

$$\underset{\underset{\sim}{c} = (c_1, \dots, c_n) \in \mathbb{R}^n}{\text{Min}} \quad \underbrace{\left\{ \sigma^2 \sum_{i=1}^{n} c_i + \lambda \left( \sum_{i=1}^{n} c_i - 1 \right) \right\}}_{\varphi_\lambda(\underset{\sim}{c})}.$$

$$\frac{\partial \varphi_\lambda(\underset{\sim}{c})}{\partial c_i} = 2\sigma^2 c_i + \lambda \quad , \quad i = 1, 2, \dots, n$$

$$\frac{\partial}{\partial \lambda} \varphi_\lambda(\underset{\sim}{c}) = \sum_{i=1}^{n} c_i - 1$$

$$\begin{cases} \dfrac{\partial}{\partial c_i} \varphi_\lambda(\underset{\sim}{c}) = 2\sigma^2 c_i + \lambda = 0 \quad , \quad i = 1, 2, \dots, n \\[2mm] \dfrac{\partial}{\partial \lambda} \varphi_\lambda(\underset{\sim}{c}) = 0 \implies \sum_{i=1}^{n} c_i = 1 \end{cases}$$

Thus $c_i = -\dfrac{\lambda}{2\sigma^2}$ , $i = 1, 2, \dots, n$ and using the last

equation $\sum_{i=1}^{n} -\dfrac{\lambda}{2\sigma^2} = 1 \implies \lambda = -\dfrac{2\sigma^2}{n}$ and

therefore $c_i = -\dfrac{\lambda}{2\sigma^2} = -\dfrac{\left(-\frac{2\sigma^2}{n}\right)}{2\sigma^2} = \dfrac{1}{n}$ , $i = 1, \dots, n$

We can further find

$$\mathcal{H} = \left[ \frac{\partial^2}{\partial c_i \partial c_j} \varphi_\lambda(\underset{\sim}{c}) \right]_{i,j = 1, \dots, n} \quad \text{and}$$

show that

$$\underset{\sim}{x}^T \mathcal{H} \underset{\sim}{x} \geq 0 \quad \text{for any } \underset{\sim}{x} \in \mathbb{R}^n$$

$$= 0 \quad \text{iff } \underset{\sim}{x} = 0$$

(14)

This then guarantees that $\underset{\sim}{c}^* = (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$ is indeed a minimizer; in fact the unique minimizer. To summarize

$$\tilde{X}_{n, \underset{\sim}{c}^*} = \sum_{i=1}^{n} \frac{1}{n} X_i = \frac{1}{n} \sum_{i=1}^{n} X_i = \bar{X}_n .$$ Thus $\bar{X}_n$ is the best unbiased linear estimator.

— Estimating Variance

So far we confined ourselves to estimation of the population mean. Now suppose we are interested in estimating variance from $X_1, \dots, X_n$ where $X_i$'s have the same mean value $\mu$, the same variance $\sigma^2$ and they are orthogonal, i.e. $\text{cov}(X_i, X_j) = 0$, $i \neq j$. A natural estimator of

$$\sigma^2 = \text{Var}(x) = E[(x - \mu)^2]$$

is its sample counterpart, i.e.

$$S_{n,*}^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X}_n)^2$$

Now first question is if $S_{n,*}^2$ is an unbiased estimator of $\sigma^2$; i.e. $E(S_{n,*}^2) = \sigma^2$

$$(X_i - \mu)^2 = [(X_i - \bar{X}_n) + (\bar{X}_n - \mu)]^2$$
$$= (X_i - \bar{X}_n)^2 + (\bar{X}_n - \mu)^2 + 2(\bar{X}_n - \mu)(X_i - \bar{X}_n)$$

$$\sum_{i=1}^{n}(X_i - \mu)^2 = \sum_{i=1}^{n}(X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu)^2 + 2(\bar{X}_n - \mu)\underbrace{\sum_{i=1}^{n}(X_i - \bar{X}_n)}_{0}$$

$$= \sum_{i=1}^{n}(X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu)^2 \quad \textcircled{I}$$

Taking expectation we find

(15)

$$E\left[\sum_{i=1}^{n}(X_i-\mu)^2\right] = E\left[n S_{n,*}^2\right] + E\left[n(\bar{X}_n-\mu)^2\right] \quad \text{(II)}$$

$$RHS = \sum_{i=1}^{n} \underbrace{E(X_i-\mu)^2}_{\sigma^2} = n\sigma^2$$

Note that $E(\bar{X}_n-\mu)=0$, i.e. $E(\bar{X}_n)=\mu$. Thus

$$E\left[n(\bar{X}_n-\mu)^2\right] = n E\left[(\bar{X}_n-\mu)^2\right] = n \, Var(\bar{X}_n).$$

On the other hand $Var(\bar{X}_n) = \dfrac{\sigma^2}{n}$. We therefore have

$$E\left[n(\bar{X}_n-\mu)^2\right] = n \cdot Var(\bar{X}_n) = n \cdot \frac{\sigma^2}{n} = \sigma^2 \text{ and hence}$$

from (II)

$$n\sigma^2 = E(n S_{n,*}^2) + \sigma^2$$

$$\Rightarrow E(S_{n,*}^2) = \left(\frac{n-1}{n}\right)\sigma^2 = \left(1-\frac{1}{n}\right)\sigma^2$$

meaning that $S_{n,*}^2$ is NOT an unbiased estimator of $\sigma$
Multiplying both sides of the last equation by the reciproc
of $(1-\frac{1}{n})$ we find $E\left(\frac{n}{n-1} S_{n,*}^2\right) = \sigma^2$. Note, howev
that $\dfrac{n}{n-1} S_{n,*}^2 = \dfrac{\not{n}}{n-1} \cdot \dfrac{1}{\not{n}} \sum_{i=1}^{n}(X_i-\bar{X}_n)^2 = \dfrac{1}{n-1}\sum_{i=1}^{n}(X_i-\bar{X}_n)^2$

Thus $\left\{ S_n^2 = \dfrac{1}{n-1}\sum_{i=1}^{n}(X_i-\bar{X}_n)^2 \right\}$ is an unbiased estimator of

Why "n-1"? "n-1" is the dimension of $\overbrace{\underbrace{\{X_i-\bar{X}_n : i=1,\cdots n\}}_{V}}^{\text{span}}$

$n-1 = \dim(span\,V)$. Note however $V$

$\dim(span\,W) = n$ where $W = \{X_i - \mu, i=1,\cdots n\}$.
We discuss these issues further in Chapter 4 where
learn regression.

— Two sample problems:

So far we only considered sampling from one population. We may have samples from two or more populations and may want to make inference about differences between the populations. Suppose for example, we want to study the difference between the average salaries of men and women,

| Men | Women |
|-----|-------|
| $X_1$ | $Y_1$ |
| $\vdots$ | $\vdots$ |
| $X_m$ | $Y_n$ |

where $X_i$'s have the common mean $\mu_X$ and $Y_j$'s the common mean $\mu_Y$. We want to estimate $\mu_X - \mu_Y$. The natural estimate is $\bar{X}_m - \bar{Y}_n$. Show that

$$E[\bar{X}_m - \bar{Y}_n] = \mu_X - \mu_Y$$

hence $\bar{X}_m - \bar{Y}_n$ is an unbiased estimator of $\mu_X - \mu_Y$. Assume further that $Xs$ and $Ys$ are independent, $Xs$ have common variance $\sigma_X^2$, $Ys$ have common variance $\sigma_Y^2$, $Cov(X_i, X_j) = 0$, $i \neq j$, $Cov(Y_i, Y_j) = 0$, $i \neq j$. Find $Var(\bar{X}_m - \bar{Y}_n)$. Hint: Use Tm 5.12.

The difference between two proportions can be treated similarly. Note that proportions are essentially means of binary variables.

(17)