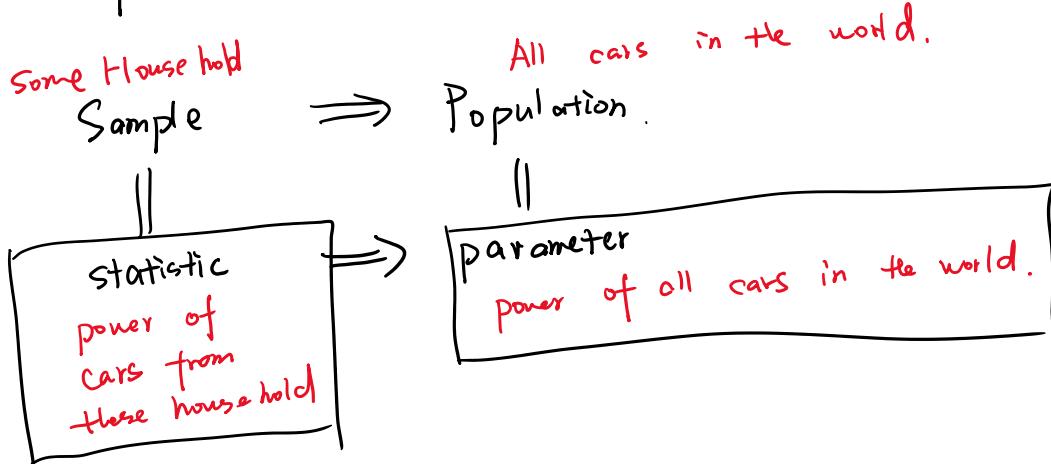


Population: The target group you want to analyse.

Sample: A small group of individuals from the population.



6.1 Numerical Summary:

Sample Mean:

Suppose we have a data set $x_1, x_2, x_3, \dots, x_n$.

The Sample mean of them is

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad x \text{ bar} \Rightarrow M$$

\uparrow statistic \uparrow parameter

Population Mean:

$$\mu = \frac{\sum_{i=1}^N x_i}{N}, \quad N \text{ is size of population.}$$

Sample Variance:

Example: 55, 63, 72, 41, 87, 75, 64, 60

$$\bar{x} = 64.625$$

$$S^2 = \frac{55^2 + 63^2 + \dots + 66^2 - 8 \cdot (64 \cdot 625)}{8 - 1}$$

$$\sigma^2 = \frac{\sum_{i=1}^{n-1} (x_i - \mu)^2}{n-1}$$

$$\sigma^2 = \frac{\sum_{i=1}^{N-1} (x_i - \mu)^2}{N}$$

s = sample standard deviation

σ = population standard deviation.

Range:

Range = largest obs - smallest obs.
 $= 87 - 41 = 46$

6.2 Stem and Leaf Plots

Example: 55, 63, 72, 41, 87, 75, 64, 60.

Median:

is the middle value of an ordered data set.

\nwarrow n is odd. $\frac{n+1}{2}$ th number

\searrow n is even average of middle two numbers.

41, 55, 60, 63, 64, 72, 75, 87

$n=8$

$$\text{median} = \frac{63 + 64}{2} = 63.5$$

Stem | Leaf

first # = 1	→ 1 4	1
first two # = 2	→ 2 5	5
last two # = 3	(3) 6	0 3 4
	→ 3 7	2 5
	1 8	7
		Last row # = 1

10 = interval length.

1 = leaf unit

$[40, 50), [50, 60) \dots \dots [70, 80)$

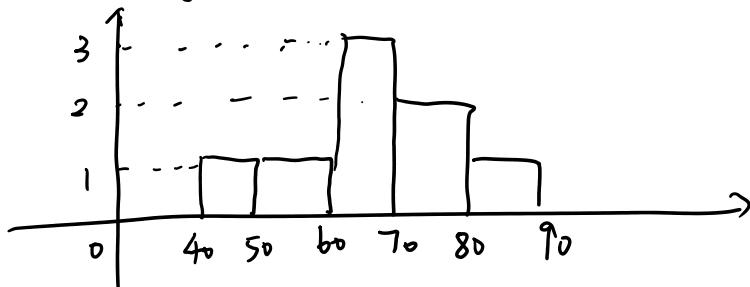
Convention when the endpoints overlap the right endpoint is not included.

6.3 Frequency Distribution and Histogram:

Interval	Frequency	Cumulative Frequency	Relative Frequency
40 - 50	1	1	$\frac{1}{8}$
50 - 60	1	2	$\frac{1}{8}$
60 - 70	3	5	$\frac{3}{8}$
70 - 80	2	7	$\frac{2}{8}$
		8	1

$80-90$	1	7	$\frac{1}{8}$
$70-80$	2	8	$\frac{1}{8}$
$60-70$	3		

Histogram.



6.4 Quartiles and Box plots

Quartile: the k th percentile P_k is the number that $k\%$ of the data is less than P_k .

median = 50% quartile, P_{50} .

Q_1 25% quartile = $\frac{n+1}{4}$ th number

Q_3 75% quartile = $\frac{3(n+1)}{4}$ th number

Example: 15, 13, 6, 5, 12, 50, 22, 18

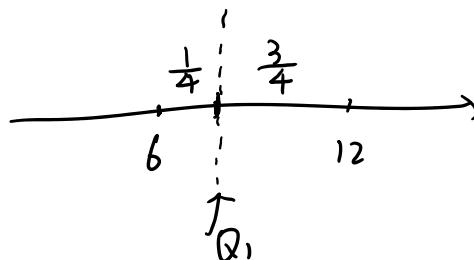
5, 6, 12, 13, 15, 18, 22, 50.

$$\text{median} = \frac{13+15}{2} = 14 \quad \text{2nd} \quad \text{3rd}$$

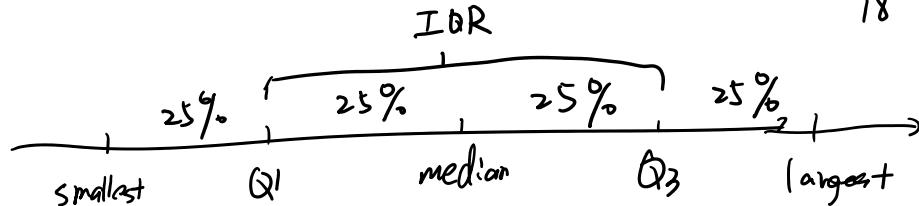
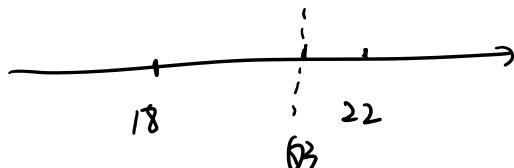
$$Q_1 : \frac{8+1}{4} = 2.25 \text{th} \Rightarrow 6 \quad 12$$

$$Q_3 : \frac{3(8+1)}{4} = 6.75 \text{th} \Rightarrow$$

$$Q_1 = 6 + \frac{1}{4} \cdot (12 - 6) = 7.5$$



$$Q_3 = 18 + \frac{3}{4}(22 - 18) = 21$$



IQR: $IQR = \text{Interquartile Range}$

$$= Q_3 - Q_1$$

Outlier: Extreme large or small value.
unusual

outlier $\left\{ \begin{array}{l} > Q_3 + 1.5 \text{ IQR} \\ < Q_1 - 1.5 \text{ IQR} \end{array} \right.$

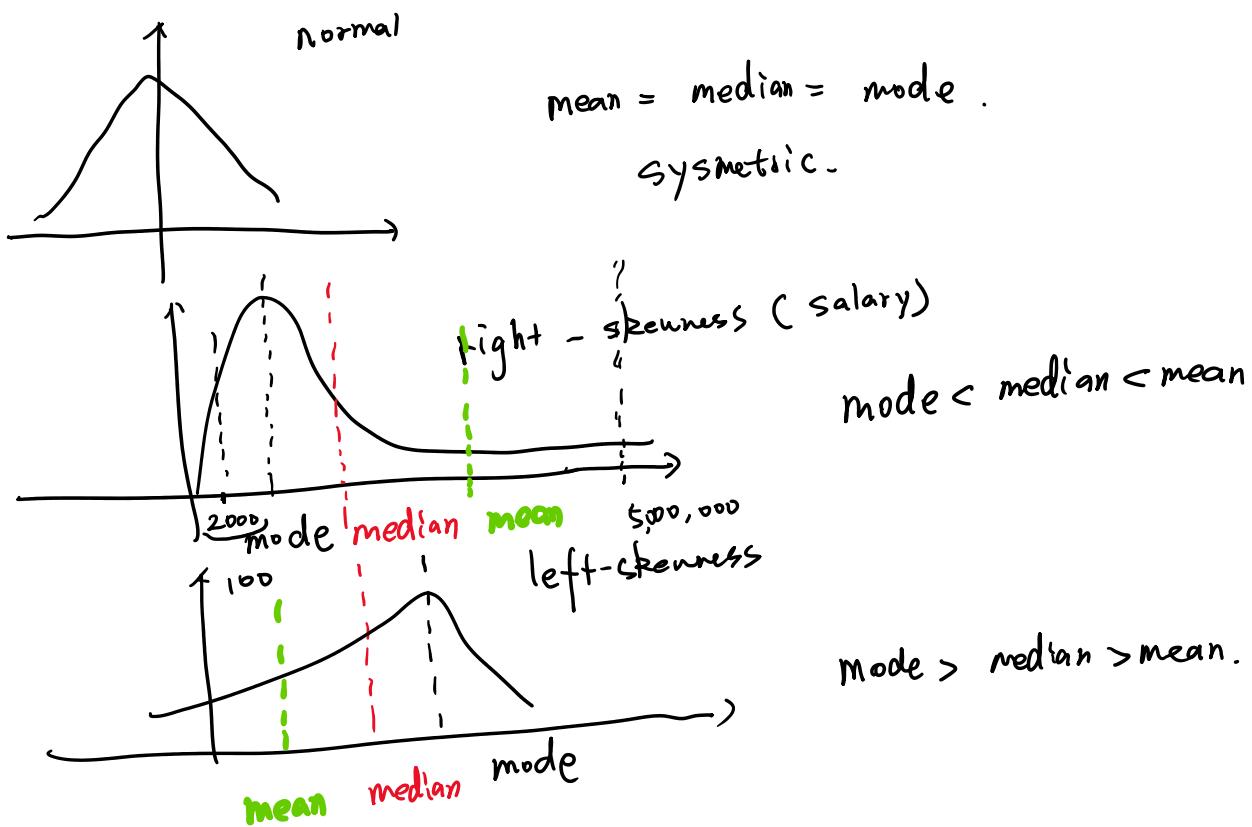
Mode: The most frequent number in the data set.

$$Q_3 + 1.5 \text{ IQR} = 21 + 1.5(21 - 7.5) = 41.25$$

$$Q_1 - 1.5 \text{ IQR} = 7.5 - 1.5(21 - 7.5) = -12.75$$

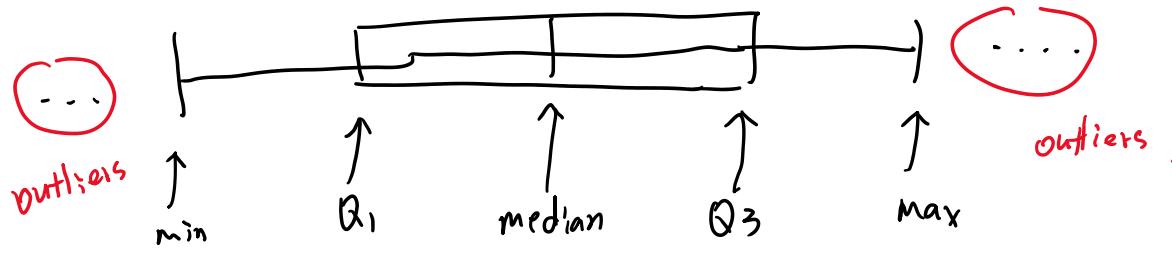
$\therefore 50 > 41.25$, 50 is an outlier

News: Average Salary of Toronto is \$50,000.



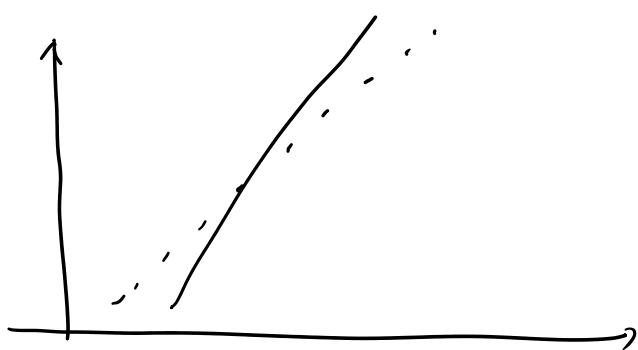
Boxplot / five number summary.

min Q1 median Q3 max.
outliers



6.7 Probability Plot.

Normal Probability Plot: Check whether sample follows Normal.



Given a sample, x_1, \dots, x_n . Find.

$$z_j = \Phi^{-1}\left(\frac{j - \frac{1}{2}}{n}\right), j = 1, 2, \dots, n.$$

$\Phi(x) = F(x)$ for a standard normal.

A normal probability plot consists n pairs $(x_{(j)}, z_j)$ where $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ are ordered data.

Example:

201, 214, 190, 176, 185, 205, 220, 183, 191, 192

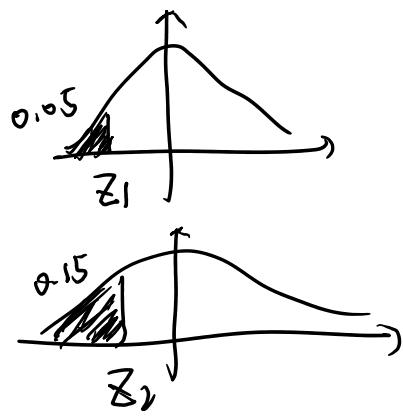
(n=10)

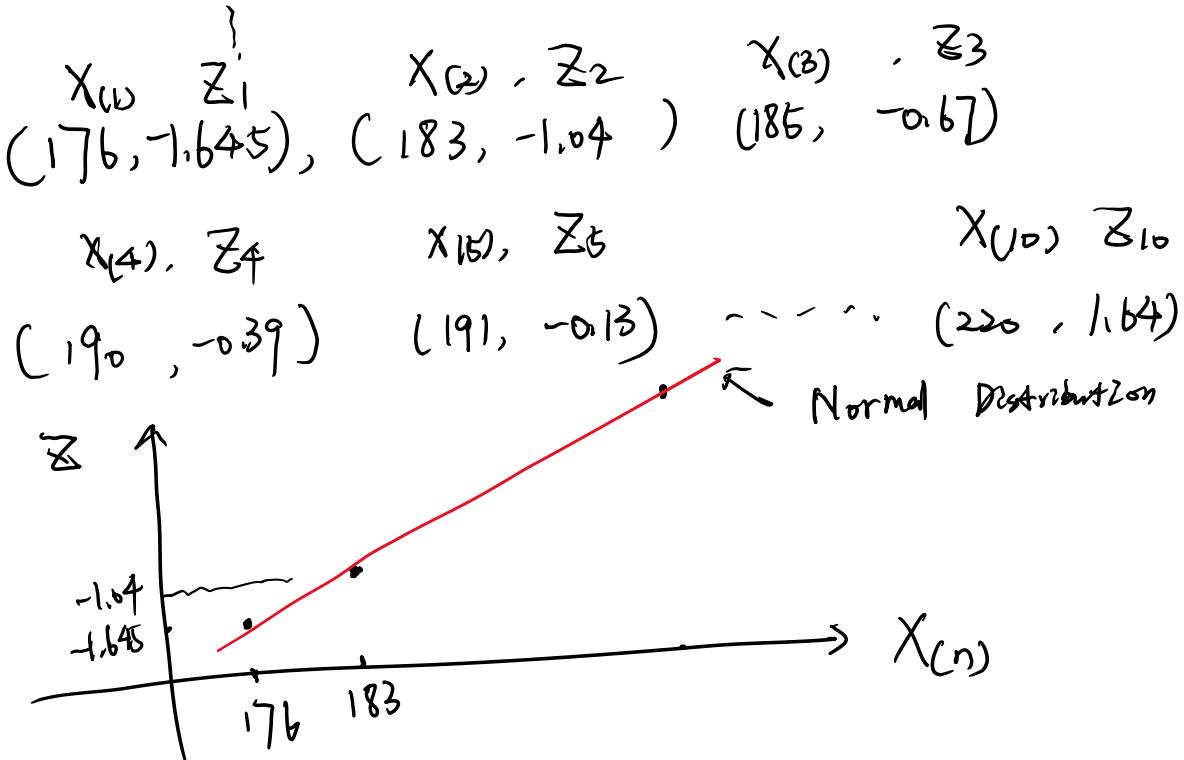
$$z_1 = \Phi^{-1}\left(\frac{1 - \frac{1}{2}}{10}\right) = \Phi^{-1}(0.05)$$

$$z_1 = -1.645$$

$$z_2 = \Phi^{-1}\left(\frac{2 - \frac{1}{2}}{10}\right) = \Phi^{-1}(0.15)$$

$$z_2 = -1.04$$





Note: interpreting the plot:

If the data comes from a population that follows normal distribution, then the points on the plot should fall close to a straight line with no discernible curve pattern. The points should be randomly scattered about the line.

7.2 Central Limit Theorem.

Definition: X_1, \dots, X_n is a random sample of size n if they are independent and all have the same distribution.

The statistic's distribution or the sampling distribution of \bar{X} follows:

$$\bar{X} \sim \text{Normal} \left(\mu, \frac{\sigma^2}{n} \right)$$

$$\text{where } E(\bar{X}) = \mu, \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

Example:

An electronic company manufactures resistors that mean μ and standard deviation 10Ω .

An electronic company manufactures resistors which have mean 100Ω and standard deviation 10Ω . Find the probability that a random sample of 25 resistors has an average of the sample is fewer than 95Ω .

$$\bar{X} \Rightarrow \mu$$

$$P(\bar{X} < 95) = P(Z < \frac{95 - 100}{10/\sqrt{25}})$$

$$= P(Z < -2.5)$$

$$= 0.006210.$$

$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$

 \downarrow
 $\text{Var}(\bar{X})$

Example:
Let x_1, x_2, \dots, x_{40} be a r.s. from the following distribution:

$$f(x) = \frac{1}{2}, \quad 4 \leq x \leq 6.$$

Find the sample mean is greater than 5.1 exactly and use normal distribution to approximate it.

① $P(\bar{x} \geq 5.1) = \int_{5.1}^6 f(x) dx$ for individual.

$= \int_{5.1}^6 \frac{1}{2} dx$ not proper.

② $E(x) = \int_4^6 x \cdot \frac{1}{2} dx = 5$

$$E(\bar{x}) = \int_4^6 x^2 \cdot \frac{1}{2} dx$$

$$\text{Var}(x) = E(x^2) - [E(x)]^2 = \frac{1}{3}.$$

$$\text{Var}(\bar{x}) = \text{Var}(x) \cdot \frac{1}{40} = \frac{1}{120}.$$

Then

$$1 \text{ Var}(\bar{X}) = \text{Var}(5 \cdot \frac{1}{\sqrt{40}}) = \frac{120}{40}.$$

Then

$$P(\bar{X} > 5.1) = P(Z > \frac{5.1 - 5}{\sqrt{1/120}}) = P(Z > 1.095) \\ = 0.1357$$

Result: If \bar{X}_1 and \bar{X}_2 are sample means of independent samples from population 1 μ_1, σ_1^2 population 2 μ_2, σ_2^2 .

Then,

$$Z = \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1)$$

$$\left(Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \right)$$

$$\left. Z = \frac{X - \mu}{\sigma} \sim N(0, 1) \right)$$

Example:

A group of bus drivers have salary as follow.
based on newspaper:

$$\text{male : } \mu_1 = 5000 \quad \sigma_1 = 40$$

$$\text{female : } \mu_2 = 5050 \quad \sigma_2 = 30$$

Then Find the probability that for a sample of $n_1 = 16$ and $n_2 = 25$, the male's average salary is 25 less than female's mean.

$$P(\bar{X}_1 - \bar{X}_2 < -25) = P(Z \leq \frac{-25 - (5000 - 5050)}{\sqrt{\frac{40^2}{16} + \frac{30^2}{25}}})$$

$$P(\bar{X}_1 - \bar{X}_2 < -2.5) = P(Z \leq \sqrt{\frac{40^2}{16} + \frac{30^2}{25}}) /$$

$$= P(Z \leq -2.14)$$

$$= 1 - 0.9838.$$

$$= 0.1162$$

7.3 Concepts of Point Estimation:

If x_1, \dots, x_n is r.s from

$$f(x) = \frac{1}{\theta} e^{-x/\theta}, \quad x > 0$$

How to find θ ?

sample	\bar{X}	S	\hat{p}	$\hat{\theta}$	estimators/statistic
population	M	σ	P	θ	true parameter

mean SD proportion

Point Estimator

Definition: A point estimator of parameter θ is sample statistic $\hat{\theta}$.

Unbiased Estimator

Definition: If $E(\hat{\theta}) = \theta$, then we call $\hat{\theta}$ an unbiased estimator.

Example: Is $S^2 = \frac{\sum_{i=1}^n x_i^2 - n\bar{x}^2}{n-1}$ an unbiased estimator of σ^2 ?

$$E(S^2) = \frac{1}{n-1} E\left(\sum_{i=1}^n x_i^2 - n\bar{x}^2\right)$$

$$= \frac{1}{n-1} [n(n+2) - 1] \rightarrow 1$$

$$\begin{aligned}
 E(S^2) &= n-1 \quad \text{as } T \\
 &= \frac{1}{n-1} \left[E\left(\sum_{i=1}^n x_i^2\right) - nE(\bar{x})^2 \right] \\
 &= \frac{1}{n-1} \left[nE(x_i^2) - nE(\bar{x})^2 \right] \\
 \text{where } &\begin{cases} E(\bar{x})^2 = \text{Var}(\bar{x}) + [E(\bar{x})]^2 = \frac{\sigma^2}{n} + \mu^2 \\ E(x^2) = \text{Var}(x) + [E(x)]^2 = \sigma^2 + \mu^2 \end{cases} \\
 &= \frac{1}{n-1} \left[n(\sigma^2 + \mu^2) - n\left[\frac{\sigma^2}{n} + \mu^2\right] \right] \\
 &= \frac{1}{n-1} [(n-1)\sigma^2] \\
 &= \sigma^2 \\
 \therefore S^2 \text{ is an unbiased estimator of } \sigma^2
 \end{aligned}$$

Compare estimators:

If two estimators are unbiased then the estimator with smaller variance is better.

$$SD(\hat{\theta}) = SE(\hat{\theta}) = \sigma_{\hat{\theta}} = \sqrt{\text{Var}(\hat{\theta})}$$

Example: Let p denote a population proportion. Then if I repeat the experiment n times and find x times success.

Then we use $\hat{p} = \frac{x}{n}$ as an estimator of p .

Find the bias and $SD(\hat{p})$.

$$\text{Bias} = E(\hat{\theta}) - \theta$$

$$\begin{aligned}
 &= E(\hat{p}) - p \\
 &= E\left(\frac{x}{n}\right) - p
 \end{aligned}$$

$(X \sim \text{Binomial}(n, p))$
 $E(X) = np$

$$\begin{aligned}
 &= E\left(\frac{\bar{x}}{n}\right) - p \\
 &= \frac{1}{n} E(x) - p \quad \leftarrow \left(\begin{array}{l} E(x) = np \\ \end{array} \right) \\
 &= \left(\frac{1}{n} \cdot np\right) - p \\
 &= 0
 \end{aligned}$$

$$SD(\hat{p}) = SE(\hat{p}) = SE\left(\frac{\bar{x}}{n}\right)$$

$$\begin{aligned}
 &= \sqrt{\frac{Var(x)}{n^2}} \\
 &= \sqrt{\frac{np(1-p)}{n^2}} \quad \left(\begin{array}{l} \text{where} \\ Var(x) = np(1-p) \end{array} \right) \\
 &= \sqrt{\frac{p(1-p)}{n}}
 \end{aligned}$$

$$\begin{aligned}
 n &\rightarrow \infty \\
 &\rightarrow 0
 \end{aligned}$$

The Mean Square Error of $\hat{\theta}$:

$$\begin{aligned}
 MSE(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\
 &= E\left[\left(\hat{\theta} - \underbrace{E(\hat{\theta})}_{\theta} + \underbrace{E(\hat{\theta}) - \theta}\right)^2\right] \\
 &= E[(\hat{\theta} - E(\hat{\theta}))^2] + E[(E(\hat{\theta}) - \theta)^2] \\
 &\quad + \underbrace{E[2(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta)]}_{0}
 \end{aligned}$$

$$= Var(\hat{\theta}) + [Bias(\hat{\theta})]^2$$

$$\Leftrightarrow E(\bar{x}^2) = Var(\bar{x}) + [E(\bar{x})]^2$$

$$\Leftrightarrow E(x^2) = \text{Var}(x) + [E(x)]^2$$

$$MSE(\hat{\theta}) = \text{Var}(\hat{\theta}) + [\text{Bias}(\hat{\theta})]^2$$

Relative Efficiency of estimators $\hat{\theta}_1, \hat{\theta}_2$ for θ

$$A = \frac{MSE(\hat{\theta}_1)}{MSE(\hat{\theta}_2)}$$

If $A > 1$, $\hat{\theta}_1$ is more than $\hat{\theta}_2$
wise versa.

(21 - 23) Sample Test 2.