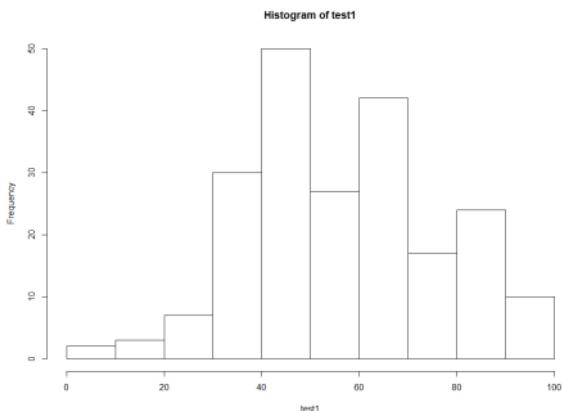


Lecture 6

May 9, 2018 6:26 PM



BSB 244 or JHE.
249.

Chapter 6: Descriptive Statistics

6.1 Numerical Summary.

Sample Mean.
n observations.
 x_1, \dots, x_n .

Sample \bar{x} \downarrow
Statistic \Rightarrow parameter μ

Population \uparrow

$$\bar{x} = \frac{x_1 + \dots + x_n}{n} = \frac{\sum x_i}{n} \Rightarrow \mu.$$

Sample variance.

$$s^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n-1}$$

$$\Rightarrow \sigma^2.$$

$$\Rightarrow \sigma^2.$$

Degree of freedom

$$S \leftarrow n-1 \text{ degree of freedom}$$
$$(x_1 - \bar{x}) + (x_2 - \bar{x}) + \dots + (x_n - \bar{x}) = 0$$
$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$
$$n\bar{x} = x_1 + \dots + x_n$$
$$(x_1 - \bar{x}) + (x_2 - \bar{x}) + \dots + (x_n - \bar{x}) = 0$$

Population mean and variance.

$$\mu = \frac{x_1 + \dots + x_N}{N}$$

$$\sigma^2 = \frac{(x_1 - \mu)^2 + \dots + (x_N - \mu)^2}{N}$$

Sample Range.

x_1, \dots, x_n . range:

$$r = \max(x_i) - \min(x_i)$$

Median:

a measure of central tendency
that divide the data into equal parts,
half below the median and half above
median.

... 2 4 . 5 . 7

median.

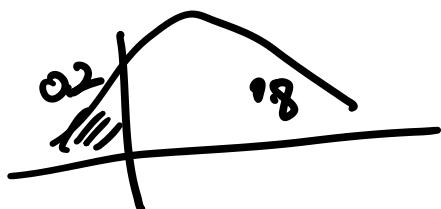
Even : 1, 2, 3, 4, 5, 6.

$$\text{median} = \frac{3+4}{2}$$

odd : 1, 2, 3, 4, 5

$$\text{median} = 3.$$

Quantile:



20th quantile
20% quantile.

$$Q_1 = \frac{n+1}{4} (\text{th.})$$

$$Q_3 = \frac{3(n+1)}{4} (\text{th.})$$

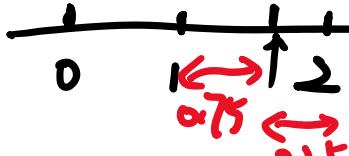
Even: 1, 2, 3, 4, 5, 6.

$$n=b \quad Q_1 = \frac{7}{4} \text{ th}$$

$$= 1.75 \text{ (th)}$$

1.75

1.75th



1 2 3 4 5 6 7

$$Q_1 = \frac{n+1}{4} = \frac{8}{4} = 2$$

$$Q_3 = \frac{3(n+1)}{4} = 6$$

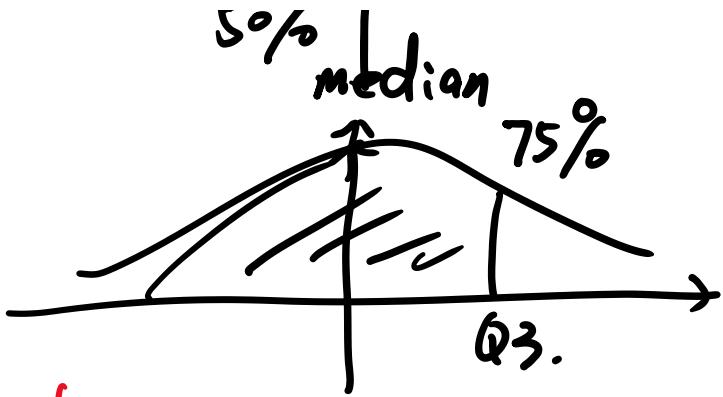
$$Q_1 = 1 \cdot 0.25 + 2 \cdot 0.75$$

Weighted.

Q₁

25% Q₁.

50% median



\downarrow 1 10 16 18 20 21.

$$Q_1 = \frac{n+1}{4} = 1.75 .$$

$\underbrace{1}_{0.75} \quad \underbrace{10}_{0.25}$

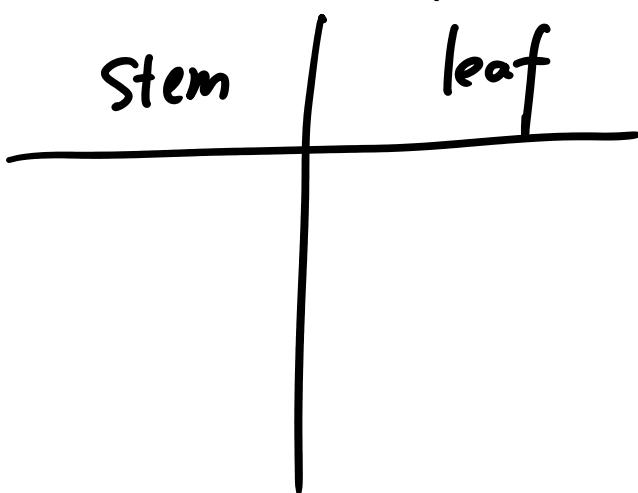
$$Q_1 = 1 \cdot 0.25 + 0.75 \cdot 10$$

Mode:

1 2 2 2 $\textcircled{3} \textcircled{3} \textcircled{3} \textcircled{3}$ 4 5 6 7

mode = 3

A stem and leaf diagram



Example:

61, 63, 64, 65, 65, 66, 70, 71

→ 1 72 75 77 78, 78, 79, 81,

71, 73, 75, 77, 78, 78, 79, 81,
 83, 84, 84, 87, 88, 88, 92, 93, 95.

Stem	Leaf	Frequency
6 6	1 3 4 5 5 6	6
(9) 7	0 1 1 3 5 7 8 8 9	9
= (7+3) 8 9 3	total sum of 2 3 5 14 4 4 7 8 8	7
3	19th 20th	3

in (the row contains median)

(last two rows.)

$$N = 25$$

$$\text{Median} = 13\text{th} = 78$$

Stem	Leaf	Fre
5 6		5
11 = 6 + 5		6
18 = 8 (5 + 6 + 7)		7
(8) 9	median	8
17 = 9 0 8		9
8		8

$$N = 43$$

$$\text{median} = 22\text{th}$$

$$Q_1 = \frac{n+1}{4} \text{ th} = 6.5 \text{ th.}$$

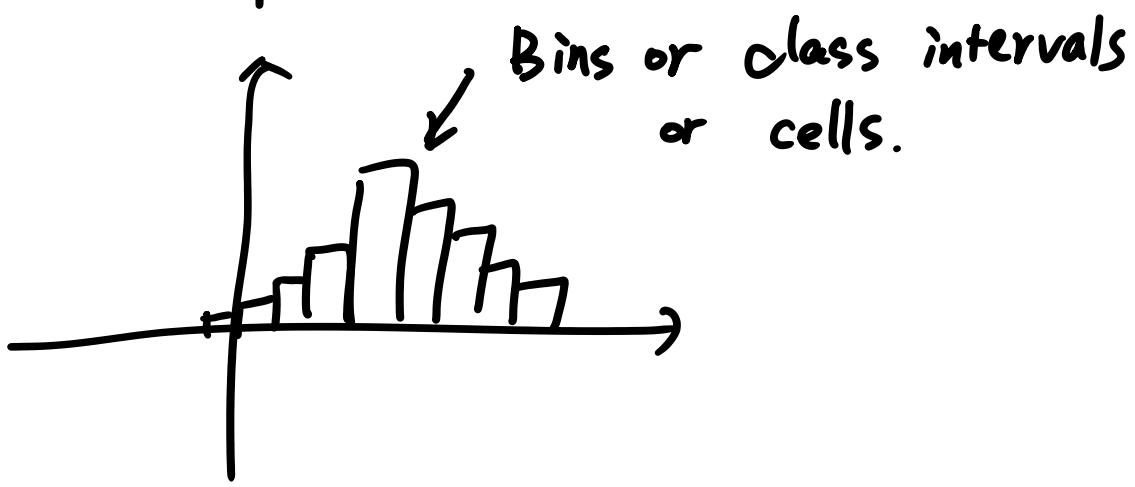
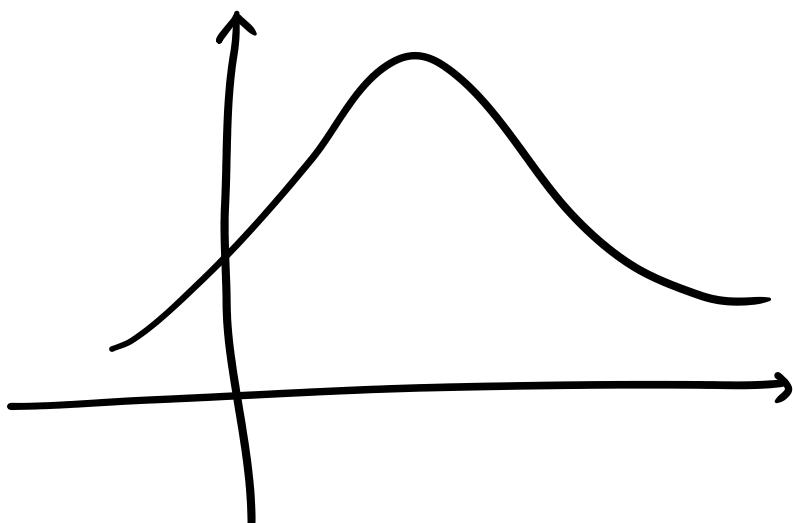
$$\frac{66 + 70}{2} = 68$$

$$Q_3 = \underline{3(n+1)} \text{ th} = 19.5 \text{ th.}$$

$$84 + 87 - \text{err}$$

$$Q_3 = \frac{3(n+1)}{4} \text{ th} = 19.5 \text{ th. } \frac{84+87}{2} = 85.5$$

Frequency distribution and histogram.



Sample size n .

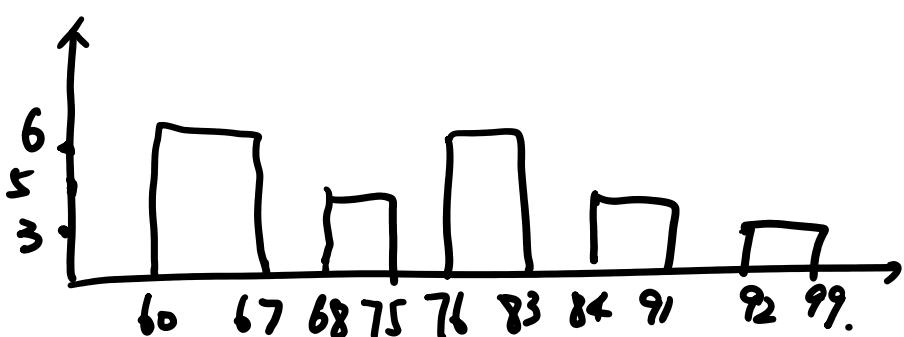
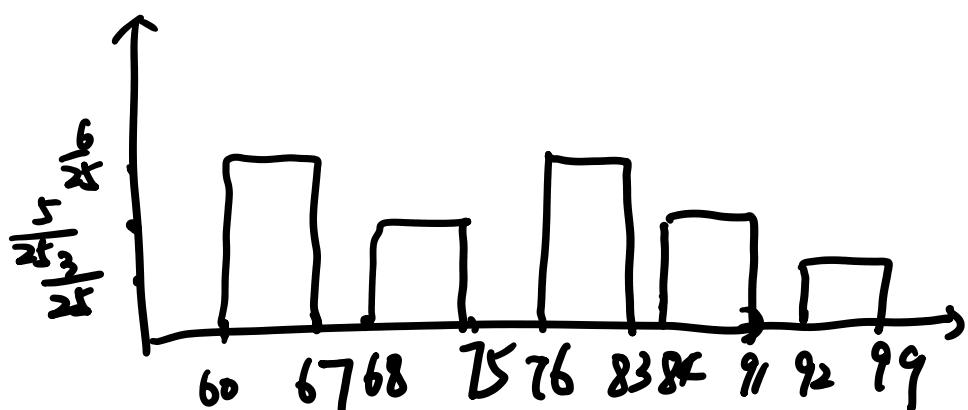
x_1, x_2, \dots, x_n .

$$\# \text{ of bins} = \sqrt{n}$$

$$\frac{\min}{61} - \frac{\max}{95}. \quad n=25$$

$$\# \text{ of bins} = 5.$$

	#	Relative Frequency
60 - 67	6	$\frac{6}{25}$
68 - 75	5	$\frac{5}{25}$
76 - 83	6	$\frac{6}{25}$
84 - 91	5	$\frac{5}{25}$
92 - 99	3	$\frac{3}{25}$



Cumulative frequency plot.

$$6 = 6.$$

$$6+5 = 11$$

$$6+5+6 = 17$$

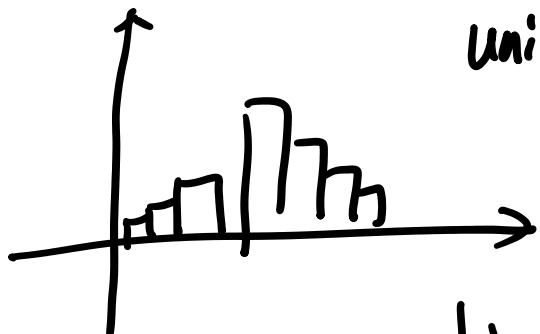
$$6+5+6+5 = 22$$

$$6+5+6+5+3 = 25$$

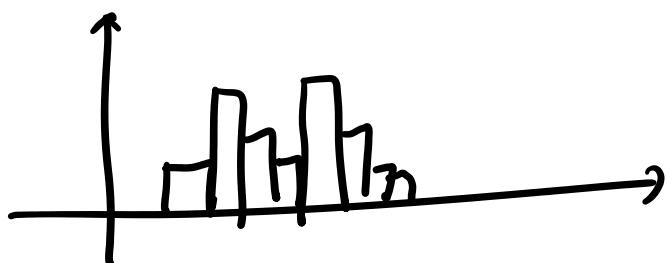
$$6 + 5 + 6 + 5 + 3 = 25$$



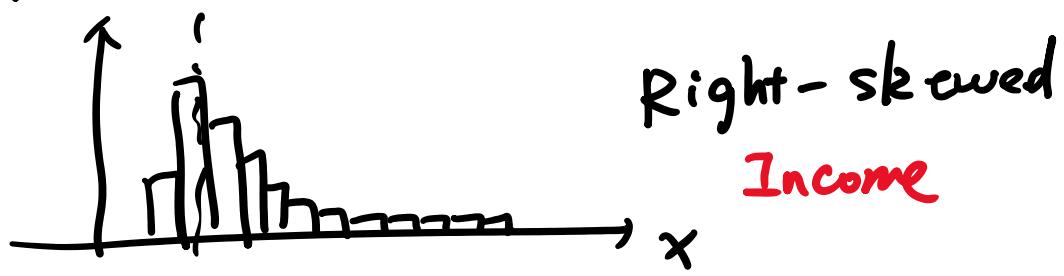
unimodal.



bimodal

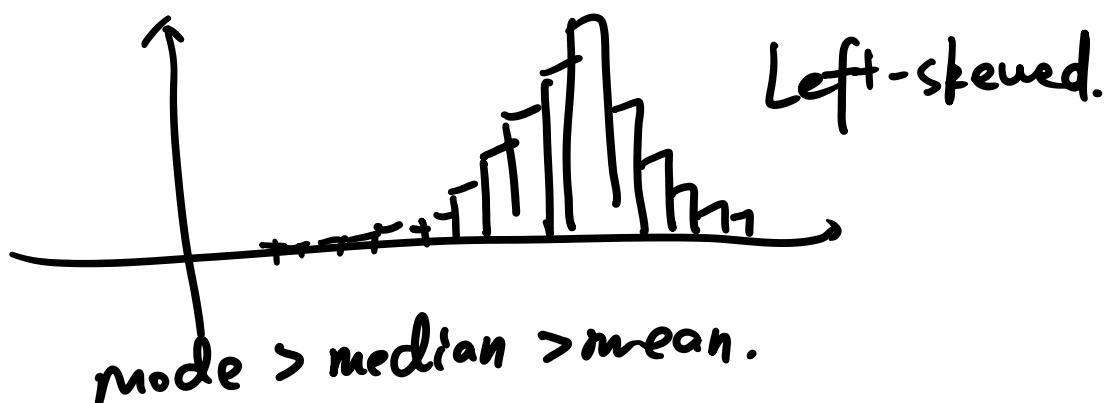


.....



mean > median > mode.

Right-skewed
Income



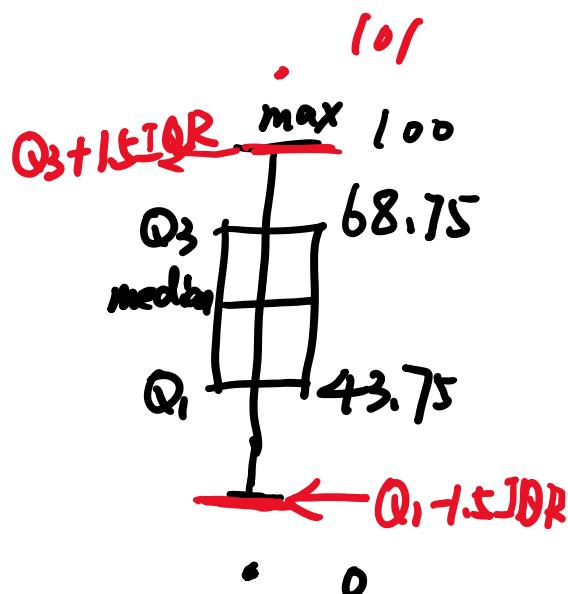
mode > median > mean.

Box plot: Box and whisker plot.

Box plot: Box and whisker plot.

$$IQR = Q_3 - Q_1$$

Interquartile Range.



If a point in the data set,

is greater than

$$Q_3 + 1.5 \text{ IQR}$$

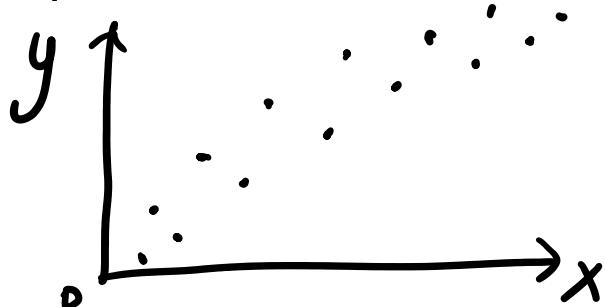
or smaller than

$$Q_1 - 1.5 \text{ IQR}$$

then it is an outlier.

Scatter plot

plot y against x .



x : weight

y : height.

Scatter plot is designed to graphically display the potential relationship between two variables.

depends on two variables.

Recall

$$r_{xy} = \frac{n\sum xy - \bar{x}\bar{y}}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

Population
 (σ_{xy})
 $\text{cov}(x,y) = S_{xy}$
 Sample

$$= \frac{S_{xy}}{S_x \cdot S_y}$$

Correlation coefficient.

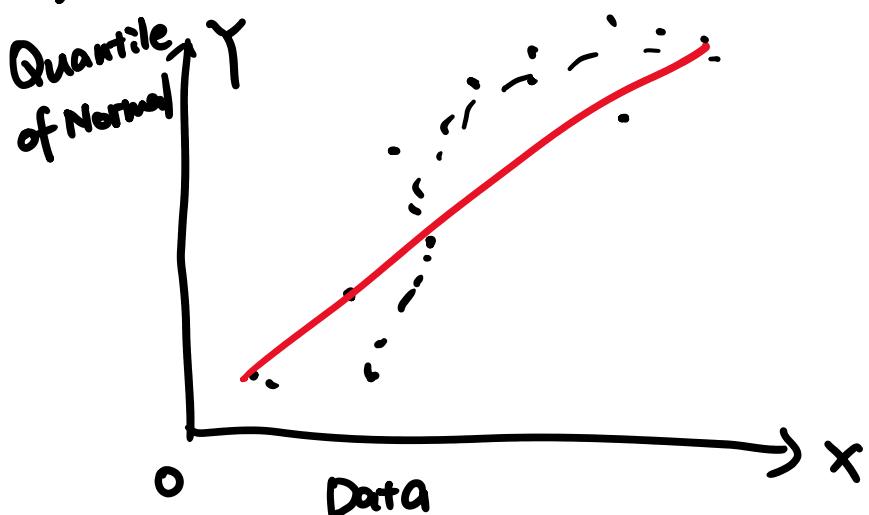
$$S_{xy} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}$$

$-1 \leq r_{xy} \leq 1$ linear association.

Normal probability Plot:

Check whether the sample follows

Normal Distribution.



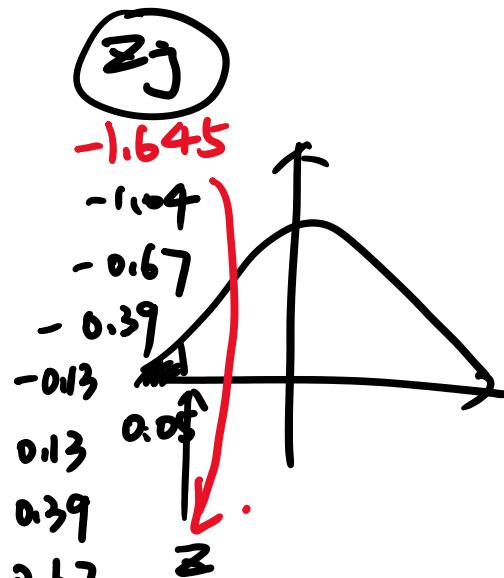
If the data set follows normal.

The Normal probability plot will falls along a straight line.

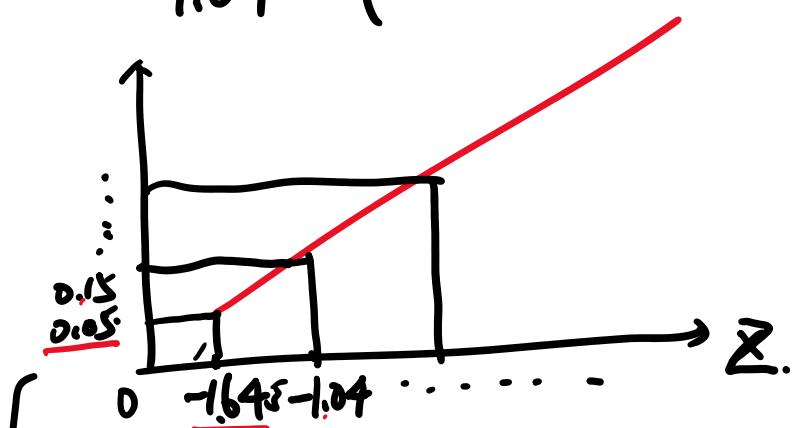
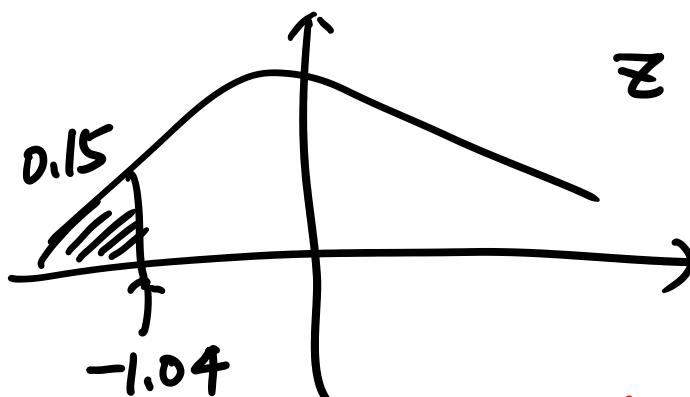
The Normal probability plot is almost on a straight line.

$$n=10$$

j	x_j	$\frac{x_j - \bar{x}}{s}$
1	176	0.105
2	183	0.15
3	185	0.25
4	190	0.35
5	191	0.45
6	192	0.55
7	201	0.65
8	205	0.75
9	214	0.85
10	220	0.95



Z - standard normal.

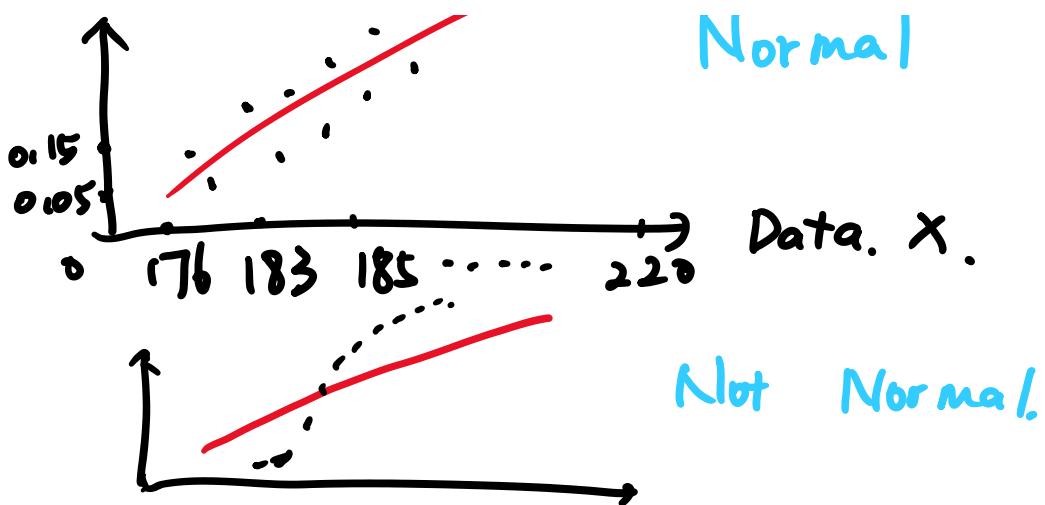


5% value smaller than -1.645

\rightarrow 5% value smaller than 5%



Normal



Chapter 7:

Point Estimation of parameters and Sampling Distribution.

$\bar{x}_1, \dots, \bar{x}_n$ Assume Poisson Distribution
then what's the λ for it.

$$\hat{\lambda} = \frac{x_1 + \dots + x_n}{n} = \bar{x} \Rightarrow \lambda$$

$\hat{\lambda}$ statistic
a function of $x_1 \dots x_n$

λ true λ .
or (true mean)
parameter

Inference

Point Estimator:

A point estimator of some population parameter Θ is a single numerical value $\hat{\Theta}$ of a statistic.

value $\hat{\theta}$ of a statistic.

Then $\hat{\theta}$ is called point estimator.

$$\hat{\mu} = \bar{x} \text{ Sample mean}$$

$\hat{s} = S$ sample standard deviation.

$$\hat{p} = \frac{x}{n} \xrightarrow{\substack{\# \text{ of } 6 \text{ for a dice.} \\ \# \text{ of total trials.}}}$$

to test whether it's a fair dice.

Sampling Distribution:

The probability distribution for a statistic is called sampling distribution.

Central Limit Theorem:

x_1, x_2, \dots, x_n are i.i.d. Sample of
 \Leftrightarrow Independent
 \Leftrightarrow identically distributed
 \Leftrightarrow Random Sample.

Size n , with mean μ and variance σ^2 .

Then.

The limiting distribution of \bar{X} .

is Normal. $(\mu, \frac{\sigma^2}{n})$.

Then.

$$\Rightarrow \frac{\bar{X} - \mu}{\sigma}$$

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

when $n \rightarrow \infty$, is a standard normal.

Examples in previous lecture note.

\bar{X}_1, \bar{X}_2 as the sample mean for two random sample.

Then, if they have σ_1^2, σ_2^2 are variance

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

as a standard normal approximately.

Example:

$$\begin{aligned} X_1 \rightarrow \mu_1 &= 5000 & \sigma_1 &= 40 & n_1 &= 16 \\ X_2 \rightarrow \mu_2 &= 5050 & \sigma_2 &= 30 & n_2 &= 25. \end{aligned}$$

Then assume X_1 and X_2 are independent. What's the probability that $\bar{X}_2 - \bar{X}_1$ is at least 25?

$$Z = \frac{(25) - (5050 - 5000)}{\sqrt{\frac{40^2}{16} + \frac{30^2}{25}}}$$

$$= -2.14$$

$$\underline{P(Z > -2.14) = 0.9838.}$$