

$X : [1, 2, 3, 4, 5]$ (Hours spent studying)

$Y : [50, 60, 70, 80, 90]$ (Exam scores)

$y_i = 50, 60, 70, 80, 90$

Calculate covariance and correlation of above.

$$COV = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{n}$$

$$COV = \frac{COV}{\sigma_X \sigma_Y}$$

x	\bar{x}	$x - \bar{x}$	y	\bar{y}	$y - \bar{y}$
1	3	-2	50	70	-20
2	3	-1	60		-10
3	3	0	70		0
4	3	1	80		10
5	3	2	90		20

$$= \frac{\sum (x - \bar{x})(y - \bar{y})}{n}$$

$$\frac{(-2)(-20) + (-1)(-10) + 0 + 1(10) + 2(20)}{5} \Rightarrow 20$$

$$\sigma_x = \sqrt{\frac{(x - \bar{x})^2}{n}} \Rightarrow \sigma_x^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

$$\text{Cov}(x, y') > \text{Cov}(x, y)$$

$\Rightarrow x$ & y' are more related.

$$r_{xy} = \frac{\text{COV}(x, y)}{\sigma_x \sigma_y}$$

$$-1 \leq r_{xy} \leq 1$$

$$= \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

$$r_{xy} = \frac{E[(x - \mu_x)(y - \mu_y)]}{\sqrt{E(x - \mu_x)^2 E(y - \mu_y)^2}}$$

$$\sum_{i=1}^n [a(x_i - \bar{x}) + (y_i - \bar{y})]^2 \geq 0$$

$$\sum_{i=1}^n \frac{a^2 (x_i - \bar{x})^2}{n} + \frac{2a}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + \sum_{i=1}^n \frac{(y_i - \bar{y})^2}{n} \geq 0$$

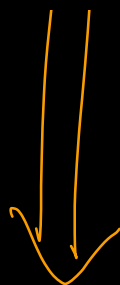
$a = \text{real no.} \neq 0$

$$a^2 \sigma_x^2 + 2a \sigma_{xy} + \sigma_y^2 \geq 0$$

$$a^2 \sigma_x^2 + 2a \sigma_{xy} + \left(\frac{\sigma_{xy}}{\sigma_x}\right)^2 + \sigma_y^2 - \left(\frac{\sigma_{xy}}{\sigma_x}\right)^2$$

$$\underbrace{\left(a \sigma_x + \frac{\sigma_{xy}}{\sigma_x}\right)^2}_0 + \underbrace{\left[\sigma_y^2 - \left(\frac{\sigma_{xy}}{\sigma_x}\right)^2\right]}_{11} \geq 0$$

$$a = -\frac{a_{ny}}{a_n}$$



$$(a_y)^2 - \left(\frac{a_{ny}}{a_n}\right)^2 \geq 0$$

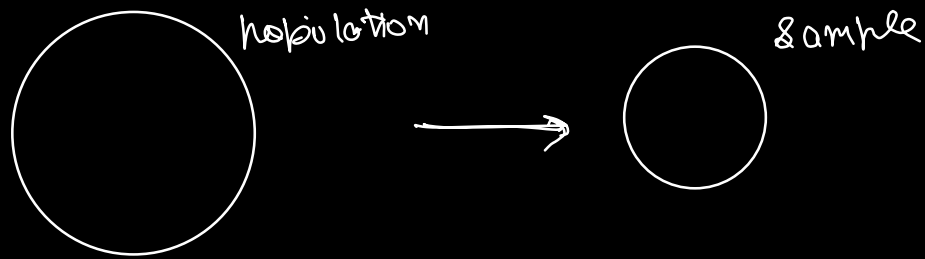
$$(a_y)^2 \geq \frac{(a_{ny})^2}{(a_n)^2}$$

$$1 \geq \frac{a_{ny}^2}{a_n^2 a_y^2}$$

$$1 \geq \gamma^2$$

$$-1 \leq \gamma \leq 1$$

Sampling Techniques



Few ways to take the sample :-

1. Simple Random Sampling :- Every member of the population (N) has a equal chance of getting selected in the sample (n).

5 stick figures $\Rightarrow 1/5$

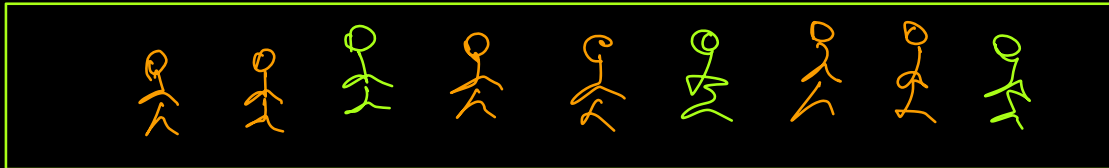
2. Systematic Sampling

We take every i^{th} individual / item

6 stick figures with red checkmarks under the first, third, and fifth figures, and red 'x' marks under the second, fourth, and sixth figures. To the right is the text $i=2$.

\times \times \checkmark \times \times \star Thanos $i=3$

Airport security check \Rightarrow 10th size



3. Stratified sampling



Strata



layers



non-overlapping region



clusters.

we will collect some proportion from each strata.

O	B	O ⁺
B ⁻	AB	AB ⁻



B⁺

[$B^+ B^+$...]

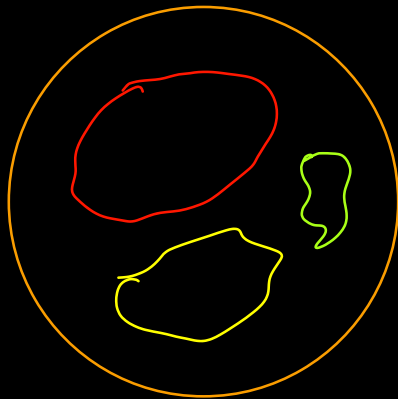
[Zero loss 0-5% 5-10% 10-20%
5% 30% ±]

3%

5%

20%

7%



✓ better sample



Convenience Sampling \Rightarrow Only those who
are interested
will be participating

apple ecosystem?

Few More Points on data

18, 20, 22, 25, 27, 28, 30, 32, 34,

35, 36, 40, 42, 45, 48, 50, 52, 55, 60, 75

80 90 100

$$\begin{aligned}\text{Range} &= \text{highest} - \text{lowest} \\ &= 75 - 18 = 57\end{aligned}$$

Percentage = per cent

represents a portion of a whole
expressed as fraction of 100

Percentile $\frac{\circ}{\circ}$

percentile actually divides
your data into 100
equal parts.

$\frac{\circ}{\circ}$ 100 $\frac{\circ}{\circ}$ ile
↓ ↓
99 marks 99% people behind you

Calculate p^{th} percentile?

$$\text{rank} = \frac{p}{100} \times \underline{(n+1)}$$

linear
interpolation

hw?

whole no. \Rightarrow take that rank

if not \Rightarrow nearest value / interpolation
as well

$$\frac{75}{100} \times (20+1) \Rightarrow \frac{21}{4} \Rightarrow \underline{5.25}$$

homework

why $n+1$

what is the 75th percentile?

Five number Summary

1. Min
2. first Quartile (25%ile) Q_1
3. Median (50%ile) Q_2
4. Third Quartile (75%ile) Q_3
5. Maximum.

IQR [Inter Quartile range]

$$= Q_3 - Q_1$$

Outlier

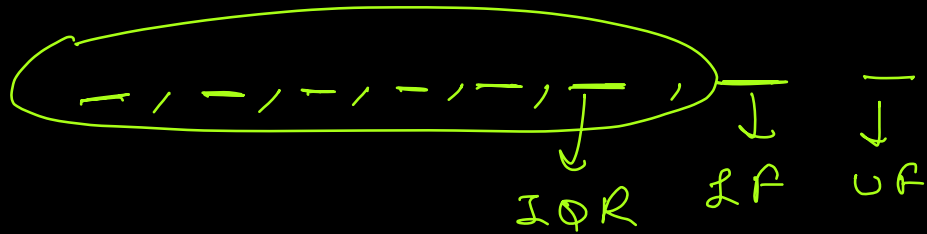


$$\text{lower fence} = Q_1 - 1.5(IQR)$$

$$\text{upper fence} = Q_3 + 1.5(IQR)$$

$$Q_1 = \frac{25}{100} \times 123$$

Q_1 Q_3
5.5, 60, 65, 70, 75, 80, 85, 90, 95, 100, 105



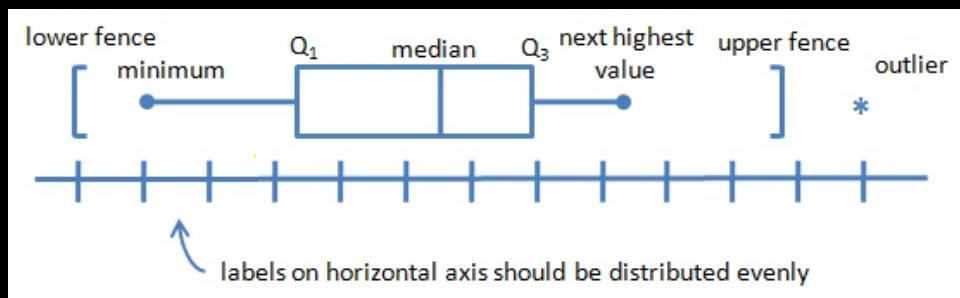
$55, 65, 80, 95, 105, 30, 20, 140$
 Q_3

$$LF = 65 - 1.5(30)$$

$$= 65 - 45 = 20$$

$$UF = 95 + 1.5(30)$$

$$= 95 + 45 \Rightarrow 140$$



Types of Graphs

- 1] Normal / Gaussian Distribution
- 2] Standard Normal Distribution
- 3] Log Normal Distribution
- 4] Power law distribution
- 5] Bernoulli's

Hypothesis →

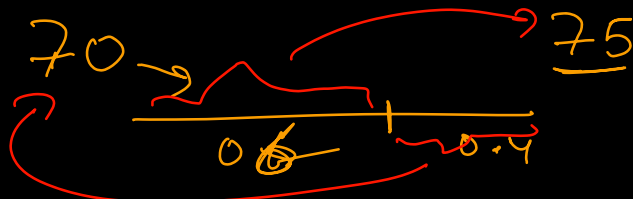
Rank Interpolation

60, 65, 70, 75, 80, 85, 90,
95, 100, 105, 110

$$\frac{30^{\text{th}} \text{ file}}{100} \times (12+1)$$

$$= 3.6$$

$$\begin{array}{c} 3.6 \\ \downarrow \\ 75.0 \approx 70 \end{array}$$



$$0.4 \times 70 + 0.6 \times 75$$

$$= 28 + 45$$

$$= 73 \Rightarrow \text{exact}$$

