

IDS Assignment- 1

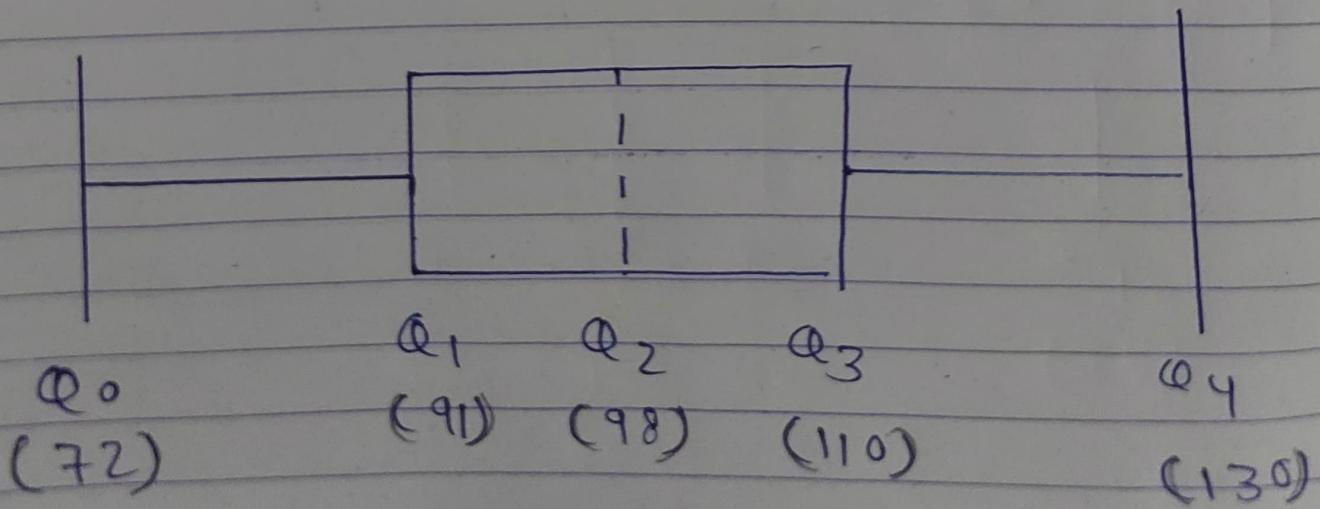
(Q1) - - - - -

SOLN

(a) - - - - -

Soln: minimum (Q_0) = 72Lower quartile (Q_1) = 91Median (Q_2) = 98Upper quartile (Q_3) = 110Maximum (Q_4) = 130

Box & Whisker Plot is:



(b) ~~Q1 = 11, Q3 = 25, IQR = 14, Median = 15~~Soln:

Sample = 25, 15, 11, 90, 14, 20, 10

In Ascending Order:

10, 11, 14, 15, 20, 25, 90

\uparrow
 Q_1

\uparrow
 Q_2

\uparrow
 Q_3

$$Q_2 = 15 \text{ (median)}$$

$$Q_1 = 11 \quad (\text{First } 25\%)$$

$$Q_3 = 25$$

$$\text{I.Q.R} = Q_3 - Q_1 = 25 - 11 = 14$$

Calculating Outliers:

$$Q_1 - 1.5 \text{ IQR} = 11 - 1.5 \times 14 = -10$$

$$Q_3 + 1.5 \text{ IQR} = 25 + 1.5 \times 14 = 46$$

$$= [-10, 46] = \text{Outliers}$$

19UCC023 - Mohit Akhouri

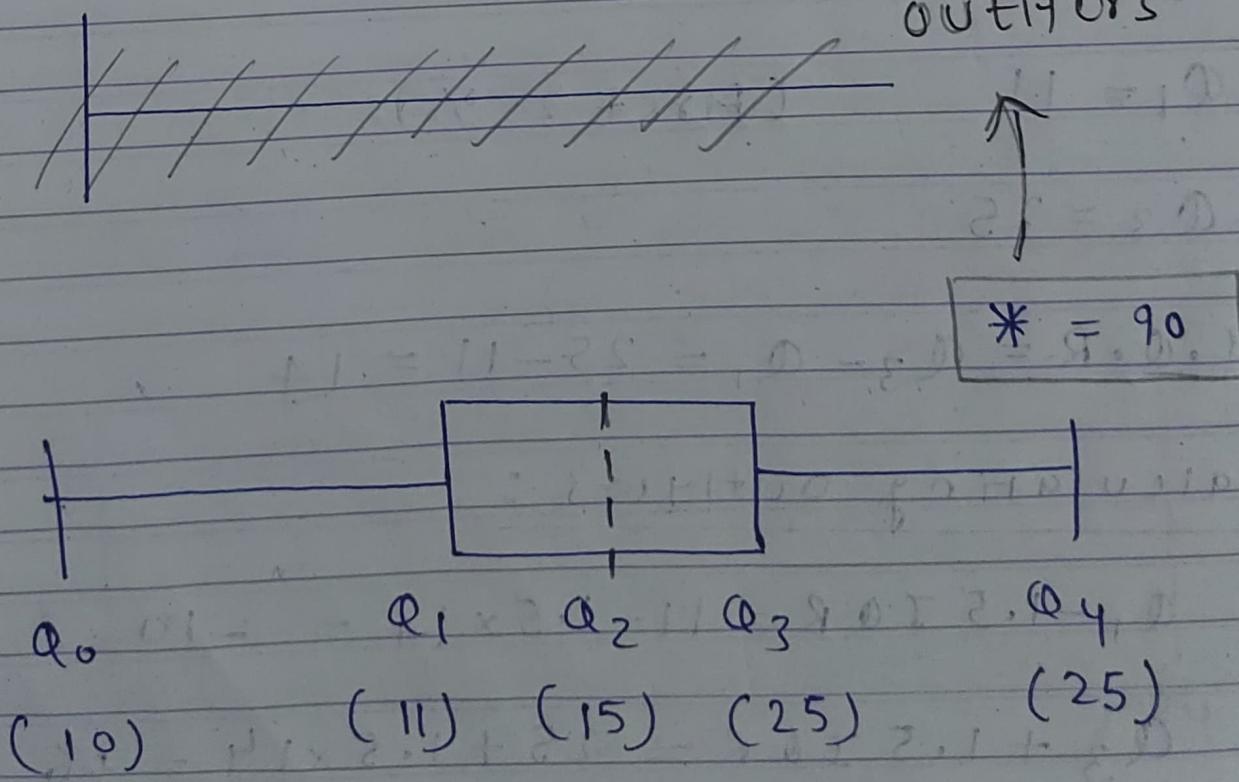
minimum (Q_0) after removing outliers

$$Q_0 = 10$$

maximum (Q_4) after removing outliers

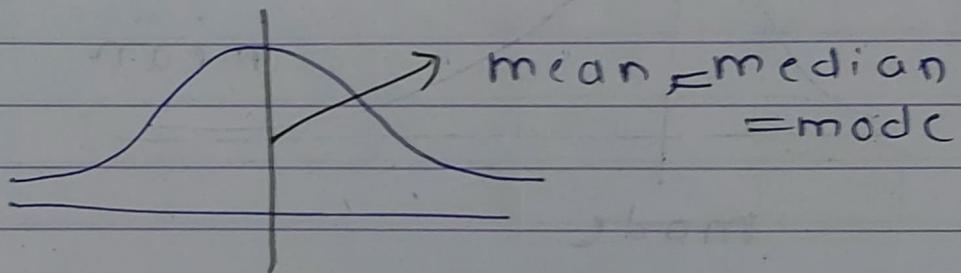
$$Q_4 = 25$$

So the box-whisker plot is:



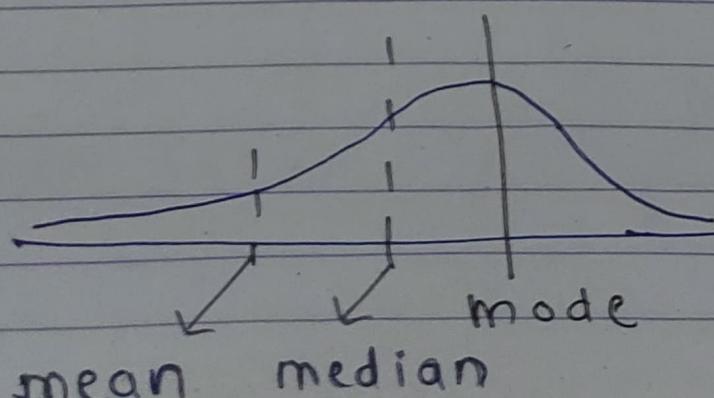
(Q2) -----

Soln: If the distribution is ~~symmetric~~ (if the distribution is ~~symmetric~~)
symmetric, all three parameters mean, median & mode are the same. (If the distribution has only 1 mode, then it is equal to mean and median, but not in case of multi-mode distribution)

LEFT-SKewed:

* If the distribution is negatively skewed (left skewed), then ~~mean~~ relationship between mean, median and mode is:

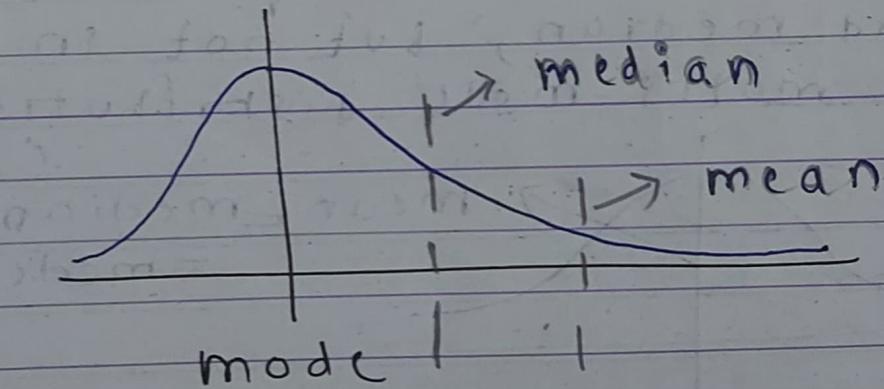
$$\text{mean} < \text{median} < \text{mode}$$



RIGHT-SKewed:

If the distribution is right-skewed (positively skewed), then relationship between mean, median and mode is:

$$\text{mode} < \text{median} < \text{mean}$$



Data pattern \rightarrow cause

Q3) -----

Soln:* Probability of head = $P(H)$

$$P(H) = \frac{40}{100} = \left(\frac{2}{5}\right)$$

* no. of times coin ~~flipped~~ flipped = (50) coin follows binomial $\sim B(n, p)$
distribution

$$\sim B(50, \frac{2}{5})$$

Now, mean of binomial random
variable is:if $X \sim B(n, p)$

then $E(X) = np$

X : tossing the coin and getting head

$$n = 50$$

$$p = \frac{2}{5}$$

$$X \sim B(n, p)$$

$$X \sim B(50, \frac{2}{5})$$

$$\boxed{\text{Here, } a = 27}$$

By Markov's Inequality,

$$P(X \geq a) \leq \frac{E(X)}{a}, \quad a > 0, \quad a \in \mathbb{R}$$

$$\Rightarrow P(X \geq 27) \leq \frac{E(X)}{27}$$

$$E(X) = np = 50 \times \frac{2}{5} = 20$$

$$\Rightarrow P(X \geq 27) \leq \frac{20}{27}$$

So, the upper bound on probability = $\frac{20}{27}$

Q4) -----

Soln:

$$\underline{n = 50}$$

$b = \frac{1}{2}$ (Since it is a fair coin)

$X = \text{coin is tossed}$

(coin follows binomial distribution)

Let $X = \text{random variable}$

~~$X = \text{coin is tossed \& we get head}$~~

$X = \text{we get head on toss of coin}$

$$X \sim B(50, \frac{1}{2})$$

We have to find probability:

$$P(X < 20 \cup X > 30)$$

$(X < 20) \cup (X > 30)$ can also be written

$$\text{as: } |X - 25| \geq 5$$

$$X \sim B(50, \frac{1}{2})$$

$$E(X) = np = 50 \times \frac{1}{2} = 25$$

$$\text{Var}(X) = np^2 = 50 \times \frac{1}{2} \times \frac{1}{2} = 12.5$$

P(X ≥ 2)

$$P(X < 20 \cup X > 30) = P(|X - 25| \geq 5)$$

So by Chebyshev's Inequality,

$$P(|X - \mu| \geq k) \leq \frac{\text{Var}(X)}{k^2}$$

$$\mu = \text{mean}, k \in \mathbb{R}, k > 0$$

$$\text{So, } P(|X - 25| \geq 5) \leq \frac{\text{Var}(X)}{5^2}$$

$$\leq \frac{12.5}{25}$$

$$P(|X - 25| \geq 5) \leq 0.5$$

So, the upper bound on probability = 0.5

Q5) -----

Soln:

$$\mu = 7$$

(a)

Soln:

~~X is student gets marks~~

$X = \text{marks of Students}$

$$(S \geq F - x_1) \hat{=} (R \geq X \geq 8)$$

~~X follows normal distribution with mean μ & variance σ^2~~

~~(approx) $X \sim N(\mu, \sigma^2)$~~ (By central limit theorem)

~~Here $\mu = 7, (S \geq F - x_1) \hat{=} P(X \geq 8)$~~

Now, we have to find probability that $X \geq 8$

$$P(X \geq 8) \leq \frac{E(X)}{8} \quad (\text{By Markov's Inequality})$$

$$E(X) = \mu = 7$$

$$P(X \geq 8) \leq 7/8$$

Upper bound on Probability = $7/8 = 0.875$

19UC023 - Mohit Akhoury

(b)

Soln:

Given $\sigma = 2, \mu = 7$

lower bound on (a)

Now, we have to find probability

$$P(5 \leq X \leq 9)$$

$$P(5 < X < 9) = P(|X - 7| \leq 2)$$

$$P(|X - 7| \geq 2) \leq \frac{\text{Var}(X)}{2^2} \quad (\text{By Chebyshev's inequality})$$

Replacing $P(|X - 7| \geq 2)$ with $P(|X - 7| \leq 2)$

$$1 - P(|X - 7| \leq 2) \leq \frac{\text{Var}(X)}{2^2}$$

$$P(|X - 7| \leq 2) \geq 1 - \frac{\text{Var}(X)}{2^2} = (8 \leq X) \geq$$

$$P(|X - 7| \leq 2) \geq 1 -$$

$$\text{Var}(X) = \sigma^2 = 4$$

$$F \geq (8 \leq X) \geq$$

So,

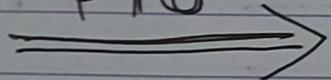
$$P(|X - 7| \leq 2) \geq 1 - \frac{4}{2^2}$$

$$P(|X - 7| \leq 2) \geq 1 - 0$$

$$P(|X - 7| \leq 2) \geq 1$$

So, the lower bound on probability is 1

PTO



Q6)

Soln:

Sample x is drawn from population with mean μ & SD σ

Mean of the sample \bar{x}

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$$\bar{x} = \frac{\mu n}{n} = \mu$$

Variance of the sample $\bar{x} = \sigma^2 / \sqrt{n}$

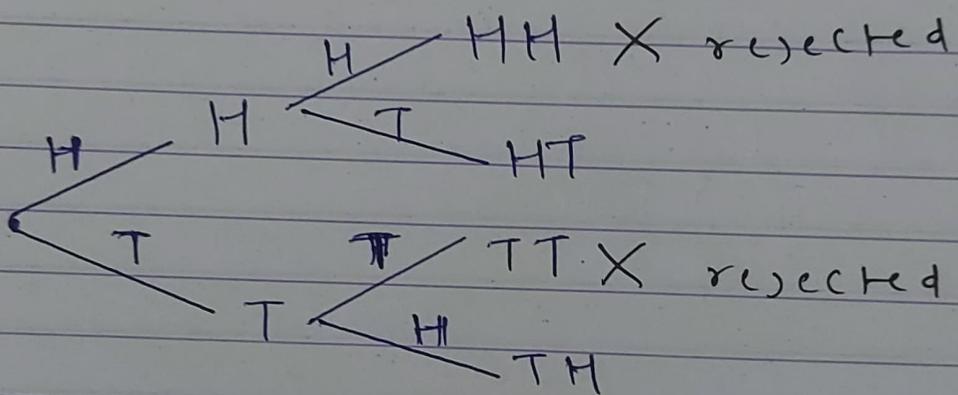
Now, if we take all such samples with mean μ & SD σ , and we combine them all, we get again a distribution which follows normal distribution due to "central limit theorem"

$$\underbrace{\quad}_{z_1} + \underbrace{\quad}_{z_2} + \dots + \underbrace{\quad}_{z_n} = \underbrace{\quad}_{N(\mu, \sigma^2)}$$

(Q7) -----

Soln: To simulate a fair coin using a biased coin, we can follow the steps given below:

(i) Toss the coin twice, we get four possible results



(ii) Now, if we get HH or TT, reject the move and try again.

(iii) Now if we get HT, TH, HT or TH, we keep (acknowledge) the first result and reject the second result.

This is how, we simulate a fair coin using a biased coin.

(Q8)

Soln: $X_1 = \text{samples of Plant 1}$ $X_2 = \text{samples of Plant 2}$ $X_3 = \text{samples of Plant 3}$ $n_1 = \text{no. of samples in } X_1 = 10$ $n_2 = \text{no. of samples in } X_2 = 10$ $n_3 = \text{no. of samples in } X_3 = 10$

$$\bar{X}_1 = \frac{\sum_{j=1}^{10} X_{1j}}{10} = 10.4$$

$$\bar{X}_2 =$$

$$\bar{X}_1 = 10.4$$

$$\bar{X}_2 = 93.3$$

$$\bar{x}_1 = \frac{\sum x_1}{n_1} = \frac{104}{10} = 10.4$$

$$\bar{x}_2 = \frac{\sum x_2}{n_2} = \frac{93.3}{10} = 9.33$$

$$\bar{x}_3 = \frac{\sum x_3}{n_3} = \frac{103.49}{10} = 10.35$$

$$\bar{x} = \frac{\bar{x}_1 + \bar{x}_2 + \bar{x}_3}{3} = \frac{30.08}{3} = 10.02$$

null hypothesis: $H_0: \mu_1 = \mu_2 = \mu_3$

$$H_0: \mu_1 = \mu_2 = \mu_3$$

$$\alpha = 0.05$$

$$SSC = \sum_{i=1}^3 (\bar{x}_i - \bar{x})^2 * n_i$$

$$\begin{aligned}
 &= (10.4 - 10.02)^2 \times 10 + (9.33 - 10.02)^2 \\
 &\quad \times 10 + (10.35 - 10.02)^2 \times 10 \\
 &= 7.294
 \end{aligned}$$

19UCC023

19UCC023 - Mohit Akhouni

$$SSE = \sum_{i=1}^3 (\text{Var}(x_i) * (n_i - 1))$$

$$= [\text{Var}(x_1) + \text{Var}(x_2) + \text{Var}(x_3)] \times g$$

$$= [3.11 + 1.68 + 1.06] \times 9$$

$$= 52.65$$

$$df_{\text{sample}} = n - 1 = 10 - 1 = 9$$

$$df_{\text{error}} = N - (k) \rightarrow \text{Total } \cancel{SSE} \text{ no. of samples}$$

$$= (10 + 10 + 10) - 3$$

$$= 27$$

$$MSC = \frac{SSC}{df_{\text{sample}}} = \frac{70.294}{9} = 0.8104$$

$$MSE = \frac{SSE}{df_{\text{error}}} = \frac{52.65}{27} = 1.95$$

$$F_{\text{stat}} = \frac{MSC}{MSE} = \frac{0.8104}{1.95} = 0.4155$$

$F_{\alpha, df \text{ col}, df \text{ within}}$

↙

↙

↙

↙

(continued from previous page)

F stat

$$F_{\text{critical}} = F_{0.05, 9, 27} = 2.2501$$

$$F_{\text{stat}} = 0.4155$$

$$F_{\text{critical}} = 2.2501$$

since $F_{\text{stat}} < F_{\text{critical}}$,

null hypothesis H_0 "failed to reject"

So, there are not enough evidences to reject the claim that all three manufacturing plants are similar.

(Q9) --- + 0.0 --- x --- 7.8 9.10

Soln:

$$n = 325$$

$$X = 225.27 \leq (n) 351.2 \text{ is not } \\ \text{above 21.02}$$

$$\sigma = 0.765$$

$$\alpha = 0.01$$

$$H - X = 5$$

$$\mu = 222$$

Step 1: Claim: $H > 222$ Opposite: $\mu \leq 222$ Step 2:

$$H_0: H = 222$$

$$H_1: H > 222$$

Step 3: $\alpha = 0.01$

Step 4:

calculate z-statistics since sample size ($n \geq 30$) and population SD is given

$$Z_{\text{stat}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

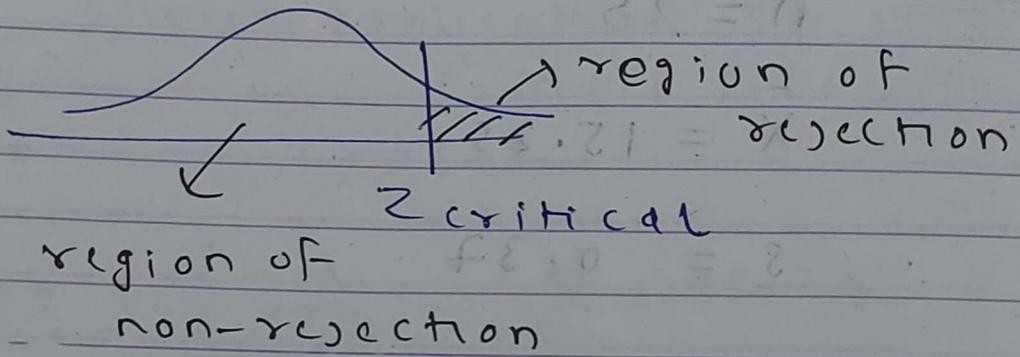
$\bar{x} = 225.27$
 $\mu = 222$
 $\sigma = 0.765$

$$= \frac{225.27 - 222}{0.765 / \sqrt{325}}$$

$$= \frac{3.27}{0.092} = 35.85$$

$$Z_{\text{stat}} = 35.85$$

Step 5: Since H_1 has ' $>$ ' sign,
go for right tailed test



$$Z_{\text{critical}} = 2.33$$

Since $Z_{\text{stat}} > Z_{\text{critical}}$,

\Rightarrow " H_0 is rejected"

\Rightarrow "alternate hypothesis failed to reject"

Step 6: Interpretation

There is enough evidences to support the claim that mean volume is greater than 222 mL.

Q10) ~~-----~~Soln%

$$\underline{n = 25}$$

$$\underline{\bar{x} = 15.25}$$

$$\underline{s = 0.37}$$

$$\underline{\mu = 15}$$

$$\underline{\alpha = 0.01}$$

Step 1:Claim: $\mu > 15$ Opposite: $\mu \leq 15$ Step 2:

$$H_0: \mu = 15$$

$$H_1: \mu > 15$$

Step 3: $\alpha = 0.01$ Step 4: Since population SD is not given and sample size < 30, go for T-test.

$$T_{\text{stat}} = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

$\bar{X} = 15.25$

$\mu = 15$

$S = 0.37$

$$= \frac{15.25 - 15}{0.37/\sqrt{25}}$$

$$= \frac{0.25}{0.074} = 3.378$$

Step 5: Since H_1 has ' $>$ ' sign, go for one-tailed test (right tailed test)

$$df = n - 1 = 25 - 1 = 24$$

$$\alpha = 0.01$$

$$T_{0.01, 24} = 2.4922 = \text{Critical}$$

one-tailed test

Step 6: $T_{\text{stat}} = 3.378$, $T_{\text{critical}} = 2.492$

Since $T_{\text{stat}} > T_{\text{critical}}$

\Rightarrow "H₀ is rejected" (reject)

$\Rightarrow H_1$ failed to reject (accepted)

Step 7: $\bar{x} = 41$

There are enough evidences to support that claim that mean weight $> 15 \text{ kg}$ is

(Q11) -----

Soln:

Different types of probabilistic sampling techniques are:

- (i) Simple random sampling
- (ii) Systematic sampling
- (iii) Stratified sampling
- (iv) Cluster sampling

(i) Simple random sampling: it can

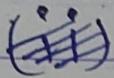
be of two types:

- (a) Sampling with replacement
- (b) Sampling without replacement

In sampling with replacement, same object may be picked multiple times but probability of picking all samples remains same.

In sampling without replacement, one sample removed from the total samples, so the probability

of choosing samples gets ^{changed} reduced by 1 each time.



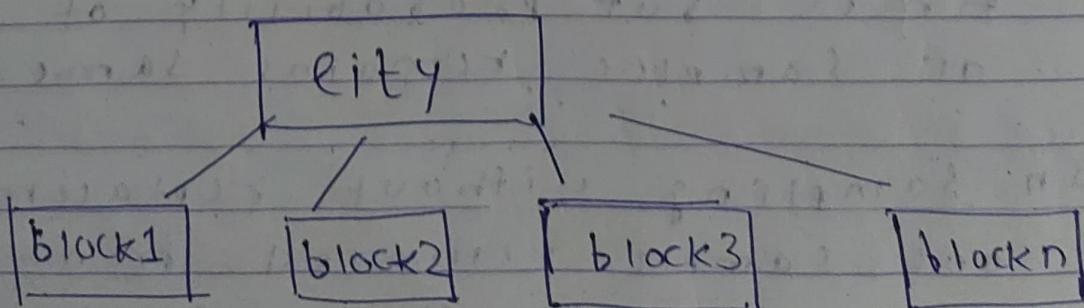
(ii) Systematic Sampling:

In this Sampling technique, we generate a uniformly distributed random number $R = U[0,1]$, then we find $i = FR \times N$. Then we choose

samples as X_{ki} where $k_i = N$ and $N = \text{size of population}$

(iii) stratified Sampling:

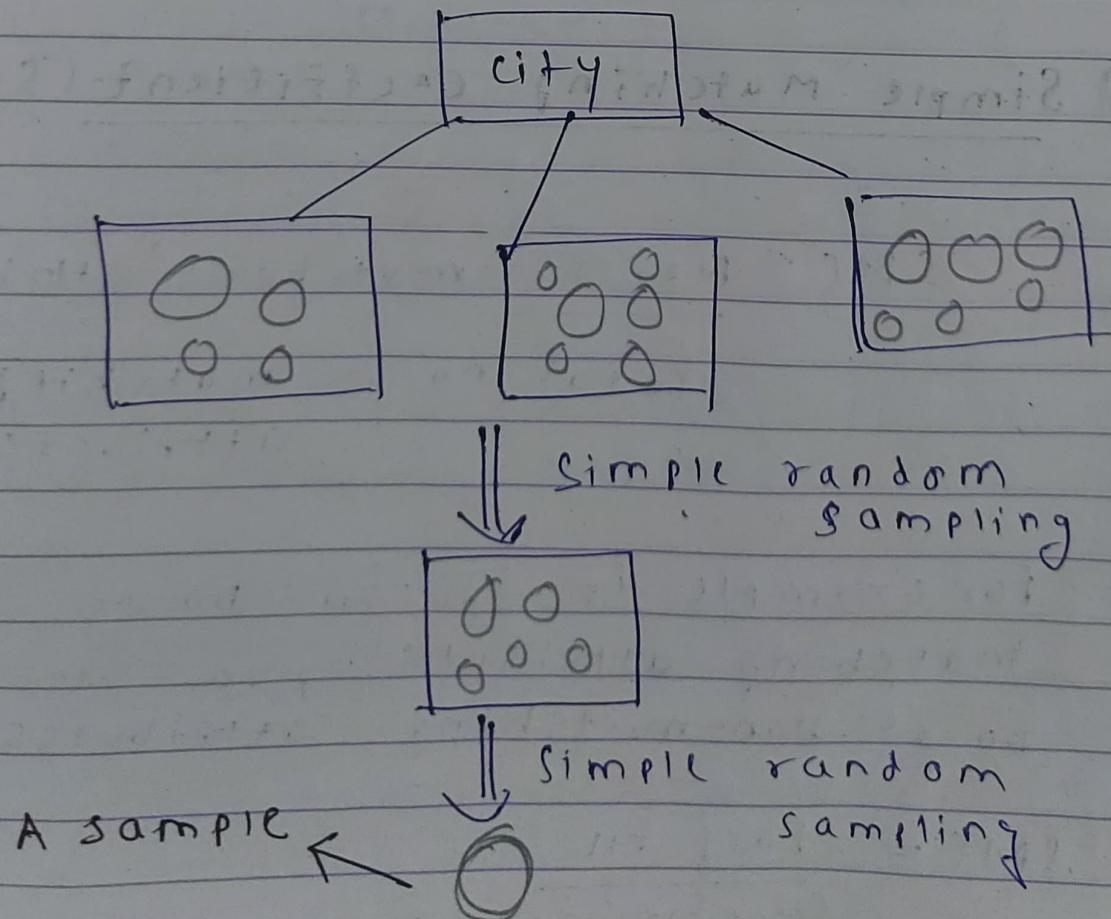
In this Sampling technique, the entire population is divided into homogenous subgroups called stratum, then we apply simple random sampling on each group.



In the above example, city is divided into blocks, then we apply simple random sampling on each block.

(iv) cluster sampling : In this sampling

technique, ~~as~~ the entire population is divided into clusters, then we select a cluster according to simple random sampling, then we apply simple random sampling within that cluster also to select samples.



(Q12) -----

SOLN:

SIMILARITY:

Similarity is defined as the measure of the degree to which two objects are alike.

similarity measure value range is between 0 and 1.

$$S \in [0, 1]$$

* MEASURES OF SIMILARITY:(i) Simple Matching Coefficient (SMC)

$$SMC = \frac{\text{No. of matching attributes}}{\text{Total no. of matching attributes}}$$

* For example if $f_{00}, f_{11} = \text{no. of matching attributes}$, $f_{01} = f_{10} = \text{no. of non-matching attributes}$

$$SMC = \frac{f_{00} + f_{11}}{f_{00} + f_{11} + f_{01} + f_{10}}$$

(ii) Jaccard Similarity Coefficient:

In this, the term f_{00} is rejected since it has no effect on the measure of similarity between two objects.

$$\text{JSC} = \frac{f_{11}}{f_{01} + f_{10} + f_{11}}$$

JSC = Jaccard Similarity coefficient

(iii) Cosine Similarity

⇒ it is a measure of similarity between 2 vectors of an inner product space.

⇒ it measures cosine of the angle between two vectors

$$\Rightarrow \cos\theta = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

DISSIMILARITY:

It is the numerical measure of degree to which two objects are different.

It also has a value between 0 and 1 and is related to the similarity coefficient by the following formula:

$$d = 1 - s \quad , d \in [0, 1]$$

Dissimilarity is measured using distance between 2 objects (feature values) which are:

$\Rightarrow r=1$ → Manhattan distance

$$\sum_{i=1}^n (x_i - y_i)$$

Minkowski
distance

$\Rightarrow r=2$ → Euclidean distance

$$\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

$$r = \left(\sum_{i=1}^n |x_i - y_i|^r \right)^{1/r}$$

$r = \text{parameter}$

$\Rightarrow r=\infty$ → supremum distance

(max. diff between

any attributes of 2 objects)

19UCC023

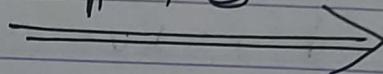
19UCC023 - Mohit Akhouri

USES OF THESE MEASURES:

These similarity & dissimilarity measures are used in various domains of data science like:

- ① TEXT-MATCHING
- ② CLUSTERING
- ③ NEAREST NEIGHBOUR

PTO



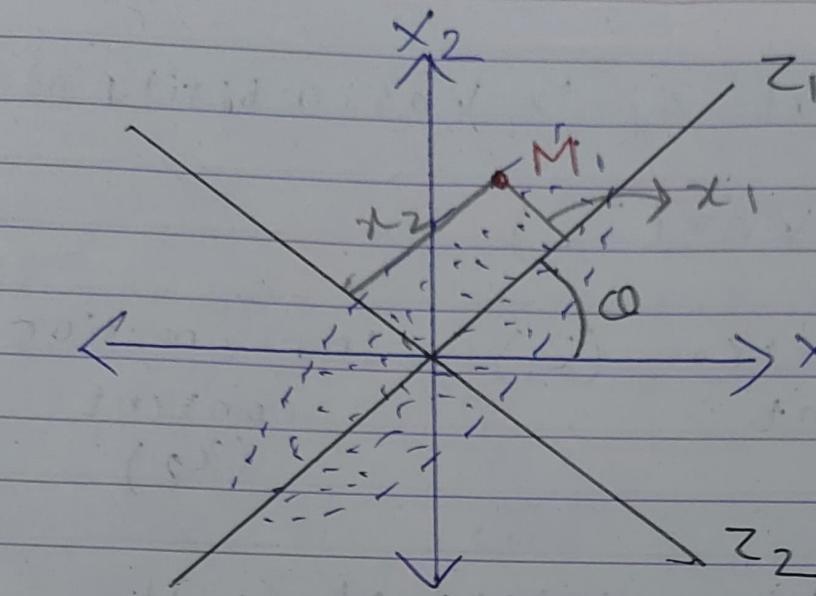
(Q13)

Soln: In dimensionality reduction, we reduce the number of dimensions to less no. of dimensions for:

- ① reducing decision taking cost
- ② increasing information content
- ③ reducing complexity of algorithm used in data science.

PCA (Principal component Analysis)

This is the most common algorithm of dimensionality reduction. In this we transform the values from x_1-x_2 plane to z_1-z_2 plane.



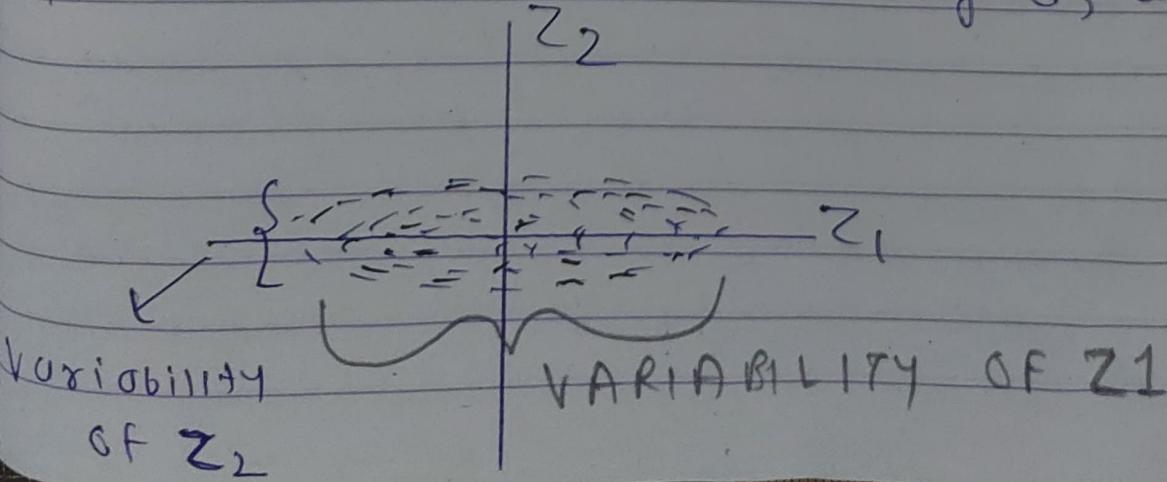
Here x_1 & x_2 are positively correlated and their projection can be found on z_1-z_2 plane as follows:

$$z_1 = x_1 \cos \theta + x_2 \sin \theta$$

$$z_2 = -x_1 \sin \theta + x_2 \cos \theta$$

PROJECTION OF POINT 'm' ON z_1-z_2 PLANE

If we rotate the axis by θ , we get:



Variability of $Z_1 >$ Variability of Z_2

So, $\text{Var}(Z_1) > \text{Var}(Z_2)$

So, information content of $Z_1 >$ information content of Z_2

Since information content of Z_1 is more, we need only variable Z_1 to take decision.

So, we reduced the dimension from 2 to 1.

This is known as "DIMENSIONALITY REDUCTION"