

"Median"

- Median depends on the position

Step 1:

- Median is the middle most value in the data set.

Steps to find the median :

Step 1: Sort data in ascending order :

10, 15, 18, 22, 25, 28, 33, 43, 50

Step 2: Find the position :

$$\therefore \frac{n+1}{2} \text{ (even or odd)}$$

$$\therefore \frac{9+1}{2} = 5^{\text{th}} \text{ order data}$$

Step 3: Find the median value .

$$\therefore 25 = \text{Median}$$

For even , 10, 15, 18, 22, 25, 28, 33, 43, 50, 1000

$$\frac{n+1}{2} = \frac{10+1}{2} = 5.5^{\text{th}} \text{ position}$$

$$\frac{25+28}{2} = 26.5$$

Quiz 1 \Rightarrow 13th November
(Monday) Wednesday

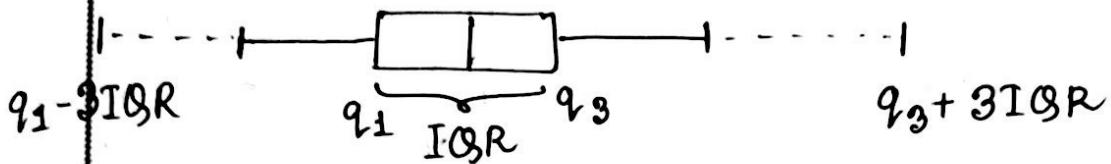
Quiz 2 \Rightarrow 2nd December
(Monday)

Quiz 3 \Rightarrow 22nd January
Wednesday

Mid \Rightarrow 19th December
Wednesday

Outlier: $\mu \pm 3\sigma$

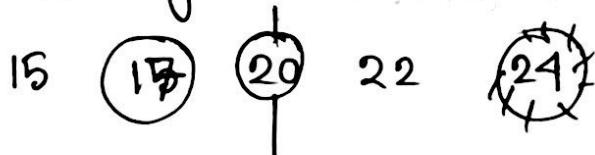
$\mu = \text{mean}$, $\sigma = \text{population standard deviation}$



Weighted Mean:- (Fixed Value का weight ध्याते हुए)

$$\bar{x}_w$$

Weighted Median:-



$$\text{Weighted median, } \bar{x} = \frac{17+20}{2} = 18.5$$

Median 20.

Mode:- (Highest)

2, 2, 1, 2, 1, 2, 1

Mode 1, 2 X

No mode ✓

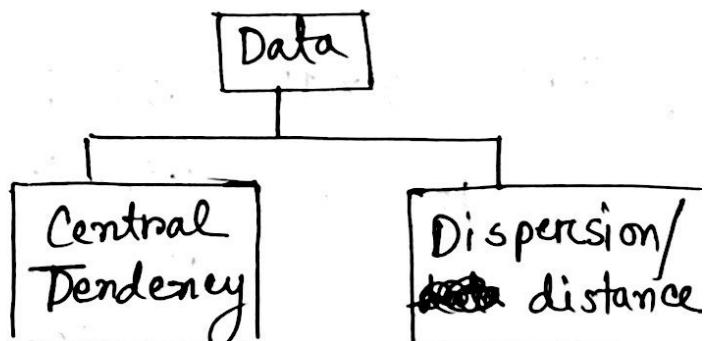
2, 2, 1, 2, 1, 2

Mode 2

Dispersion:-

Distance

can't be negative



why $(n-1)$?

Variance :- $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$

Sample :-

Practise examples

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

$$s = \sqrt{s^2} = \sqrt{\text{Variance}}$$

↓
standard deviation

Degree's of freedom :

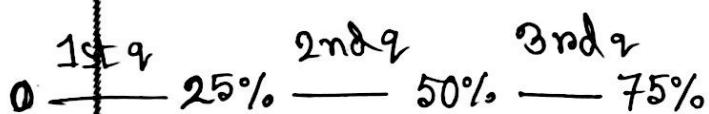
Mean absolute deviation: $\frac{\sum_{i=1}^n |x_i - \bar{x}|}{n}$

Median absolute deviation :

∴ sort → Median, → median Absolute deviation → sort → median

Practise examples homework

quantile



$$1\text{st quantile} = \frac{1}{4} = 25\%$$

position: $0.25 \times 15 = 3.75$ th value = 10

$$0.75 \times 15 = \text{th value} =$$

$$10 + 0.75(4\text{th} - 3\text{rd}) = 10 + 0.75(10 - 10)$$

position: $0.75 \times 15 = 11.25$ th value = 10

$$= 15 + 0.25(12\text{th} - 11\text{th}) = 15 + 0.25(16 - 15)$$

$$= 15 + 0.25(1)$$

$$= 15.25$$

95% percentile

$$\text{position: } 95 \times 15 = 14.25\text{th}$$

$$= 14\text{th value} + 0.25(15\text{th} - 14\text{th}) \\ = 17 + 0.25(20 - 17) = 17.75$$

Interquartile Range (IQR)

$Q_3 = 75\text{th percentile}$

$$\text{IQR} = Q_3 - Q_1$$

$$\text{Range} = \text{max} - \text{min}$$

$140-160$ $160-180$ $180-200$	\nearrow upper limit \nwarrow excluded
$50-55$ $55-60$ $61-65$	Exclusive Interval Inclusive Interval

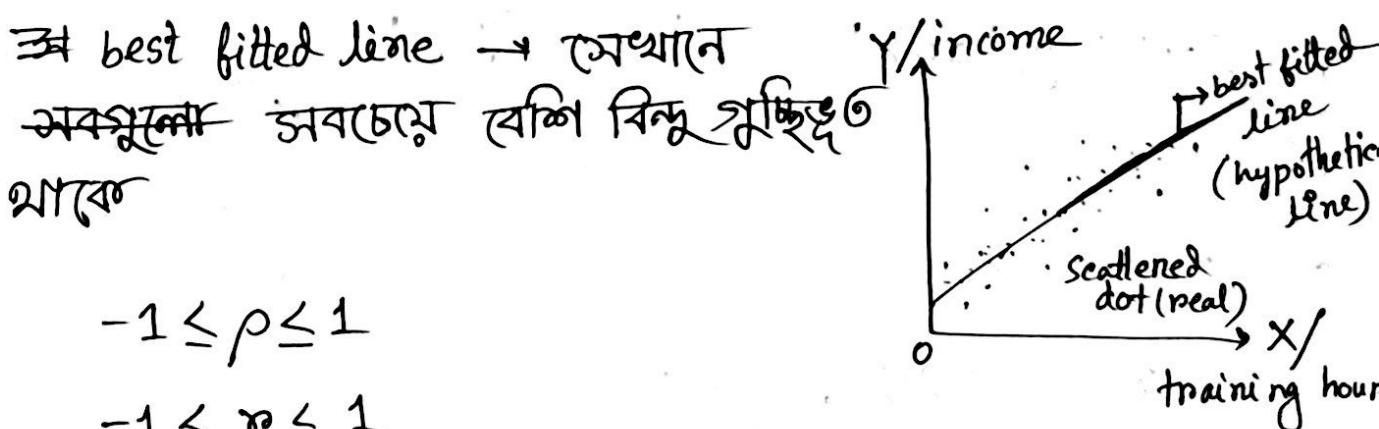
C.W

05/11/2024



Correlation

দুটি numericable variable এর অন্তর্ভুক্ত relationship.



$$-1 \leq \rho \leq 1$$

$$-1 \leq r \leq 1$$

-1 টির বেশিরভাটি গেল Strong negative linear relationship
 1 " " " Strong positive relationship

+0.5 ++

Rule of thumb

$0.5 \rightarrow$ Strong positive \nearrow linear correlation

< -0.5 Strong negative \nwarrow linear correlation

0.0002 There is no linear correlation

0.23 positive \nearrow linear correlation

~~Probability range~~ $0 \leq \text{probability range} \leq 1$

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(N-1) S_x S_y}$$

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1) S_x S_y}$$

$$S_x = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} ; \quad S_y = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n-1}}$$

Sl.no	x_i	y_i	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
1	-	-	-	-	2
2	-	-	-	-	2

* n এর মান
কখনও 0
হবে না

$r = 0.529$, There is a strong linear positive correlation.



Covariance

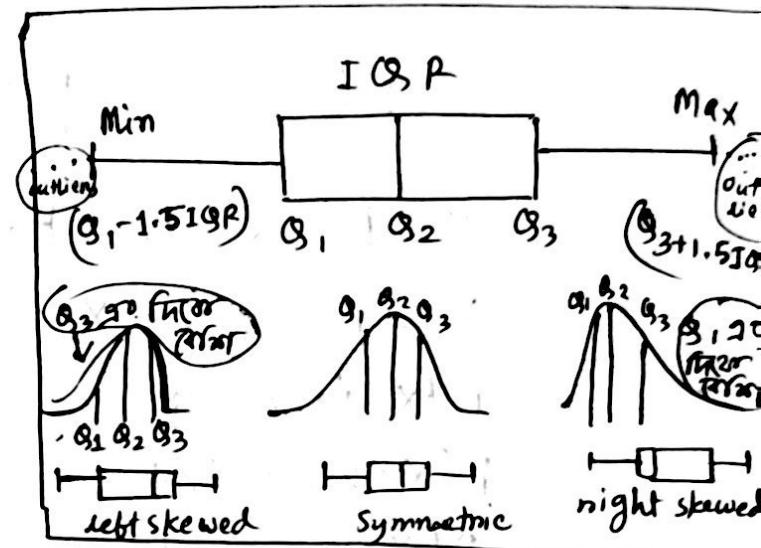
$$\sigma_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N} = S_{xx}$$

$$S_{xx} = \frac{1}{N} \sum (x_i - \bar{x})(x_i - \bar{x})$$

$$\therefore \sigma_{xy} = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)(y_i - \mu)$$

$$\therefore \rho_{xy} = \frac{\text{cov}(xy)}{\sigma_x \sigma_y}$$

$$\therefore \text{cov}(xy) = \rho \sigma_x \sigma_y$$



Lecture 3

Exploratory data analysis

Attribute: Data points or samples are described by attribute.

- ① Nominal / Categorical
- ② Ordinal
- ③ Binary
- ④ Numerical

quiz 01: Slide 1 and Slide 2

13th November 2024

Book : Hayter

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\mu = E(x) = \sum x p(x) = \sum \frac{1}{n} x = \sum p(x)x$$

$$2^2 = 4, S = \{HH, HT, TH, TT\}$$

X = no. of H, $X = 2, 1, 1, 0$

$$P(X=1) = 2/4$$

$$P(X=2) = 1/4$$

$$P(X=0) = 1/4$$

Bar graph ← Discrete Vs Continuous attribute

→ Box plot, Histogram
→ Chronology

Attribute is the data point.

$$3 \leq X \leq 10$$

$$X = -3, -2.5, 10, 12, 11.25$$

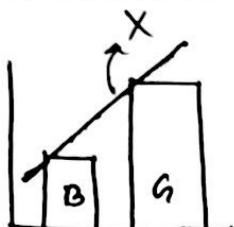
$$-2.5 \leq X \leq 12$$

Discrete can be negative

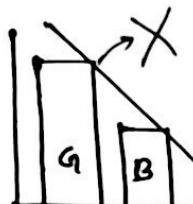
Discrete only equality matters. But continuous doesn't matter.

Continuous $\{P(X > 2) = P(X > 2)\}$

Discrete $\left\{ \begin{array}{l} 1 - P(X < 2) \\ \Rightarrow 1 - [P(X=0) + P(X=1)] \end{array} \right.$



Bar graph



Bar graph



Histogram

H (Effect)

	O(No)	Yes(1)		
No(0)	50	25	75	
S(Cause)	Yes(1)	20	55	75
	70	80		

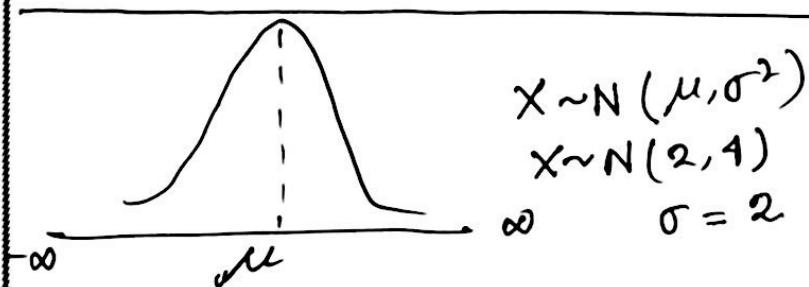
X

Data distribution

Normal distribution :- / Gaussian distribution

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

Discrete : pmf : probability mass function



Normal Distribution
 * mean = median = mode
 * symmetry

Continuous :
pdf : prob density function

$$0 < \sigma^2 < \infty$$

$$Z = \frac{x-\mu}{\sigma}$$

$Z \sim N(0, 1)$

$$\text{Expectation } (a) = a$$

Median,

$$F(x) = P(x \leq w)$$

$$\text{Cdf } F(x) = 0.5$$

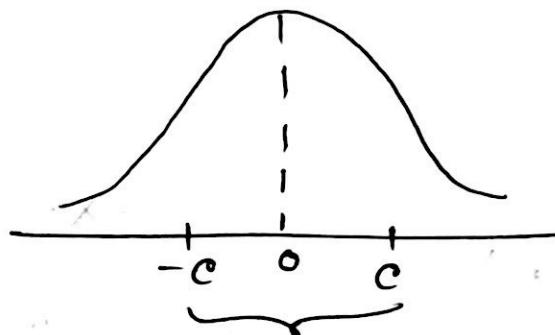
cdf up to x

18/11/2024
Properties

Normal Distribution

Example-1

$$\mu \pm \frac{\sigma z_{\alpha/2}}{c}$$



$$95\% = 2\sigma$$

$$68\% = 1\sigma$$

$$99.7\% = 3\sigma$$

$$X \sim N(\mu, \sigma^2)$$

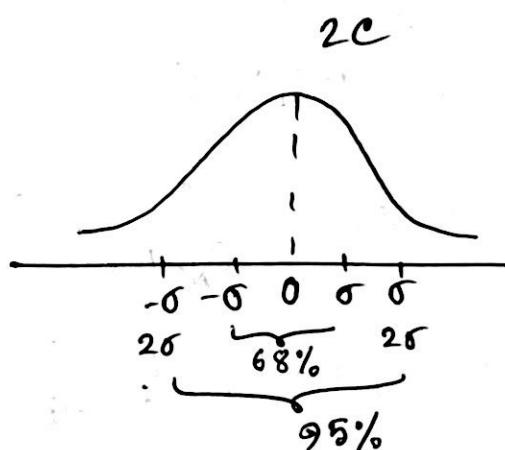
$$Z \sim N(0, 1)$$

σ or value 1

$$Z = \frac{X - \mu}{\sigma}$$

$$X \sim N(2, 4)$$

~~$$Z \sim N(1.4, 0.15^2)$$~~



$$Q(z) = P(Z \leq z)$$

$$F(x) = P(X \leq x)$$

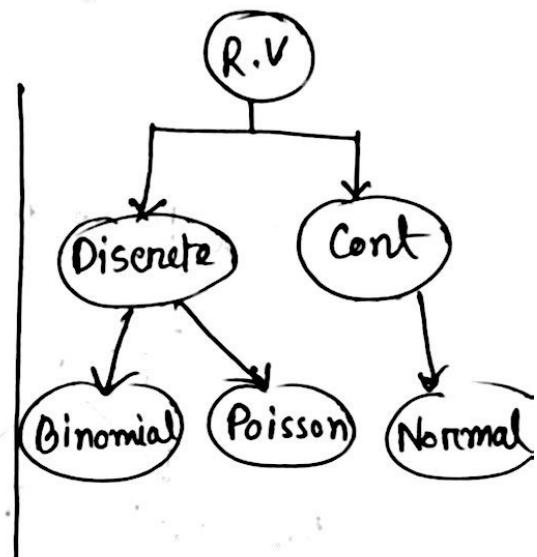
$$50\% = F(x)$$

$$\begin{aligned}\hat{\mu} &= \bar{x} \\ \hat{\sigma} &= s\end{aligned}\quad \left. \right\} \text{Estimated}$$

$$\bar{x} = \frac{\sum x_i}{n}$$

$$s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

Cdf = Cumulative distribution
of



Homework

Binomial Distribution

$$P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$$

C.W
20/12/2024

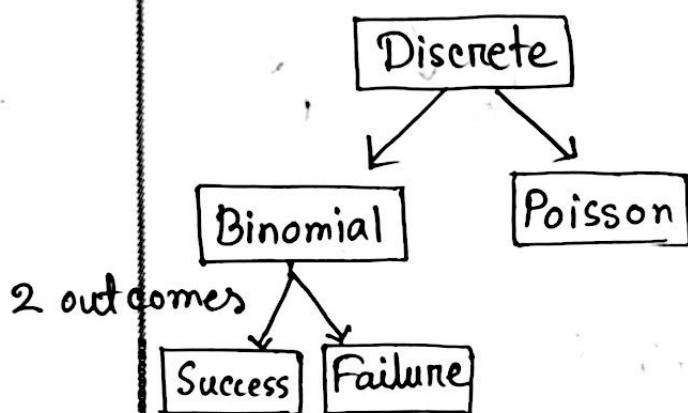
Paper Source
Subject
Date Time

Mid - 11 / 12 / 2024

Binomial Distribution

$$x = 0, 1, 2, \dots, n$$

$${}^n C_x p^x (1-p)^{n-x} ;$$



Bino: x = number of success

p = probability " "

n = no. of trial / sample size

$p(x)$ = probability mass . (pmf)

discrete এবং pmf always

$$E(X) = np$$

$$\text{Var}(X) = npq$$

$$q = 1-p$$

; q = probability of failure

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

~ follows

Poisson / Counting Distribution

Poisson

$X = \frac{\text{number of events / incident occurs during any time interval}}{\text{exponential}}$

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}; x \geq 0 \quad x = 0, 1, 2, \dots$$

Poisson limit 0 to ∞ .

$$E(X) = \text{Var}(X) = \lambda$$

λ = average no. of events occur during time interval

Practise Problem:



$$n = 10$$

$$p = 0.5$$

$$\begin{aligned} P(X=6) &= \binom{10}{6} p^6 (1-p)^{10-6} \\ &= \binom{10}{6} (0.5)^6 (0.5)^4 = 0.205 \end{aligned}$$



$$p = 0.6$$

$$n = 10$$

$$\begin{aligned} P(X=7) &= \binom{10}{7} (0.6)^7 (0.4)^3 \\ &= 0.215 \end{aligned}$$

*** for discrete, equality is everything. But for continuous, equality is nothing.

$$P(X \geq 7) = P(X = 7) + P(X = 8) + P(X = 9) + P(X = 10)$$

$$\Rightarrow 1 - P(X < 7) = 1 - [P(X = 0) + \dots + P(X = 6)]$$

Poisson: $P(x: \mu) = (e^{-\mu} * \mu^x) / x!$

⊕ $\lambda = 2$ per year

$$P(X = 3) = \frac{e^{-2} 2^3}{3!} = 0.180$$

*** Probability can't be negative, also can't be greater than 1.

⊕ Hayter Book, Chapter-3

3.1 exercise (Page - 159)

a) $P(X = 2)$

b) $P(X \geq 1) = 1 - P(X < 1)$

$P(X \geq 2)$

Answers

H.W: 3.1.4, 3.1.5, 3.1.6

Poisson: 3.4.6, 3.4.8

CORRELATION ANALYSIS (NORMAL Data)

$$\chi^2 = \sum_{\text{chi}^2}$$

$$0 \leq \chi^2 \leq \infty$$

$$E_{11} = \frac{450 \times 300}{1500} =$$

$$E_{12} = \frac{450 \times 1200}{1500} =$$

$$E_{21} = \frac{1050 \times 300}{1500} =$$

$$\chi^2 = \sum \left(\frac{(O_{ij} - E_{ij})^2}{E_{ij}} \right)$$

$$\chi^2 = \left(\frac{(250-90)^2}{90} + \frac{(200-360)^2}{360} + \frac{(50-210)^2}{210} + \frac{(1000-840)^2}{840} \right)$$

$$\chi^2_{\text{cal}} = 507.93.$$

$$df (\text{degree's of freedom}) = (r-1)(c-1) = (2-1)(2-1) = 1$$

Null Hypothesis

H_0 : No association

H_a/H_1 : There is an association.

2 numeric value association \rightarrow Correlation
 2 categorical value " $\rightarrow \chi^2$
 Cardinal
 Categorical effect

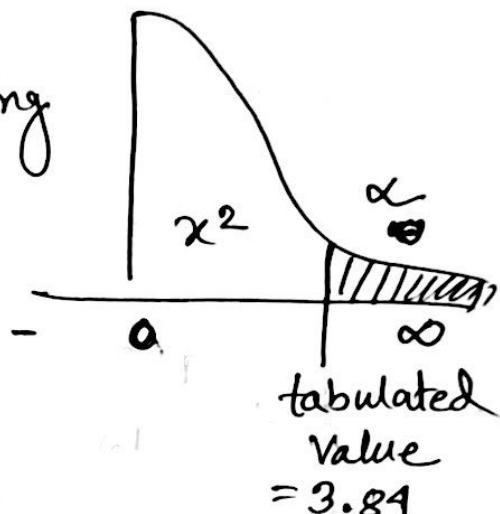
		H	H	
		R1	R2	G Total
Cause	S	250 0ii	200 012	450
	S	50 021	1000 022	1050 R2
		300	1200	1500
	C1	C2		G Total

$$O_{ij} = E_{ij} = \frac{R_i C_j}{G T = n}$$

$$E_{11} =$$

যদি α না দেওয়া থাকে, by default $\alpha = 5\%$.

H_0 : There is no association between chess playing and liking the science fiction



$$\chi_{\text{cal}}^2 > \chi_{\text{tab}}^2 \rightarrow H_0 \text{ rejected}$$

We may reject H_0 .

We have sufficient evidence that statistically and significantly, there is an association between the chess playing and liking the science fiction. Researcher wants to prove the alternative (H_A)

H.W Practise Problem

$$CI + \alpha = 1$$

→ Formula.

Hayten Book

Practice Problem

Paper Source

Subject

Date: 27/11/24 Time:

	Republican	Democrat	Independent	Total
Male	100	70	30	200
Female	110	60	20	220
Total	210	130	50	420

∴ Expected Values:-

$$\begin{array}{l|l|l} E_{11} = \frac{200 \times 210}{420} = 114.28 & E_{12} = \frac{130 \times 200}{420} = 61.90 & E_{13} = \frac{200 \times 50}{420} = 23.80 \\ E_{21} = \frac{220 \times 210}{420} = 125.71 & E_{22} = \frac{130 \times 220}{420} = 68.09 & E_{23} = \frac{220 \times 50}{420} = 26.19 \end{array}$$

∴ Calculation of χ^2 :

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$\begin{aligned} \chi^2 = & \frac{(100 - 114.28)^2}{114.28} + \frac{(70 - 61.90)^2}{61.90} + \frac{(30 - 23.80)^2}{23.80} + \frac{(110 - 125.71)^2}{125.71} \\ & + \frac{(60 - 68.09)^2}{68.09} + \frac{(20 - 26.19)^2}{26.19} \end{aligned}$$

$$\chi^2 = 4.159 + 4.0486 = 8.5076$$

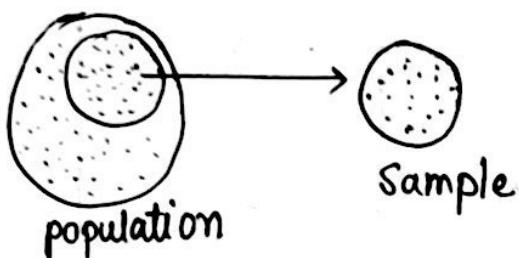
calculated

$$\text{degrees of freedom} = (2-1)(3-1) = 2$$

$$\chi^2_{\text{tabular}} = 7.824$$

$$\text{Significance level } \alpha = 0.02$$

∴ We may reject H_0 .
 There is significant relation between gender and political party. Researcher wants to prove H_0 .



$$\begin{matrix} \mu \\ \bar{x}, S, S^2 \\ \sigma \\ \rho \\ n \end{matrix}$$

$$X \sim N(\mu, \sigma^2)$$

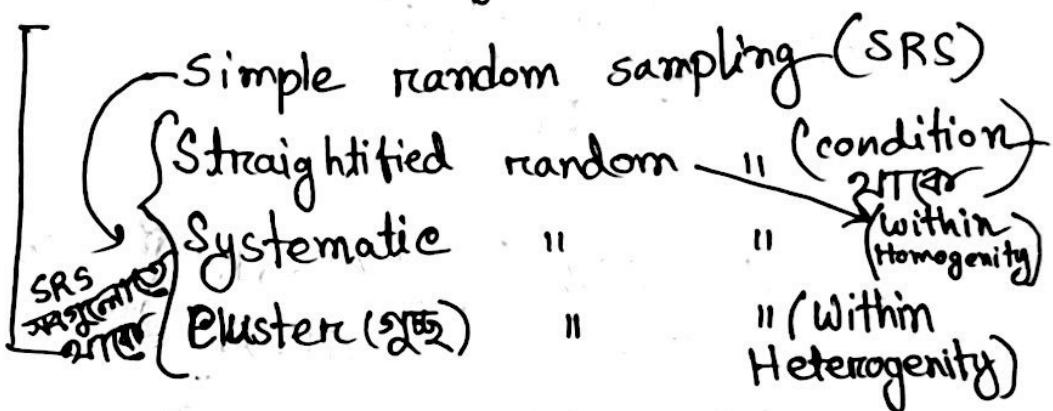
$$X \sim B(n, p)$$

$$X \sim P(d)$$

BHHS 2022
 HIES 2022
 MICS 2022

population list \rightarrow Sampling f

4 pillars of Sampling



SRS has no conditions.

$$\bar{x} = \hat{\mu}$$

$$E(\bar{x}) = \mu$$

$$E(\bar{x}) - \mu = 0 \text{ (Unbiased)}$$

$$E(\bar{x}) - \mu \neq 0 \text{ (Biased)}$$

same

$$E(\hat{\mu}) - \mu = E(\hat{\theta}) - \theta = \text{Bias} \neq 0$$

বইয়ে দেওয়া

0.0001 হলেও biased

Exhausted = Fullfilled

Mutually Exclusive = No common thing (\cap) $A \cap B = \emptyset$

Collectively Exhausted = Fullfilled (\cup) $A \cup B = S$

Book : Cochran (Sampling Test)

$$n = \left(\frac{z\sigma}{E} \right)^2$$

$E =$ margin of error

σ = standard deviation

Quiz - Slide 3 & 4. (Monday)

Mid - 11 Dec 2024

০.৫ এর নিচের value এর জন্য
negative ztable . ০.৫ এর
উপর হালে (+) z table.

Normally closest lower
value নিবৰ্ত্ত

$$\alpha = 10\% = 0.1$$

$$P = -1.65$$

$$n = \left(\frac{z\sigma}{E}\right)^2 \quad \left\{ \begin{array}{l} E = \text{margin of Error} \\ \sigma = \text{std.} \end{array} \right.$$

$$CI + \alpha = 1$$

$$CI = 90\% \Rightarrow \alpha = 10\%$$

$$n = \frac{z^2 pq}{E^2} \quad \left\{ \begin{array}{l} E = \text{margin of Error} \\ z = \text{standard normal} \\ q = 1-p \end{array} \right.$$

Sampling Distribution of a statistics :-

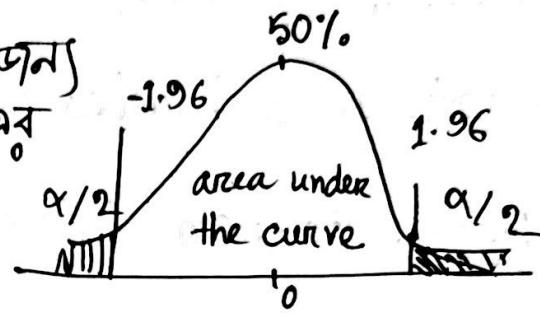
If you increase the sample size , central

$$\bar{X}, S, S^2$$

Sampling distribution of
Variance \propto square root of গুরুত্ব

$$\text{Var}(x) = SD(x)$$

$$\text{Var}(\bar{x}) = SD(\bar{x})$$



$$\alpha = 5\% = 0.05$$

$$\frac{\alpha}{2} = 0.025$$

- ① Vs code , Conda
- ② How to integrate conda with Vs code
- ③ What is python (or conda) environment ?
- ④ How to maintain multiple environment ?
- ⑤ Create an account on GitHub.
- ⑥ Install git in your device .
- ⑦ How to connect with a remote repo using Single VS code ?

Git basics : repo, branch, stage , commit, push , pull , merge

download m.
import
clone

Install necessary extension on VS code for version control (github)

Create your ~~first~~ first project and push it on github

remote device

local device → (origin)



Data analysis from Team Data set.

No classes on Monday.

Sampling distribution of a statistics

Sample size

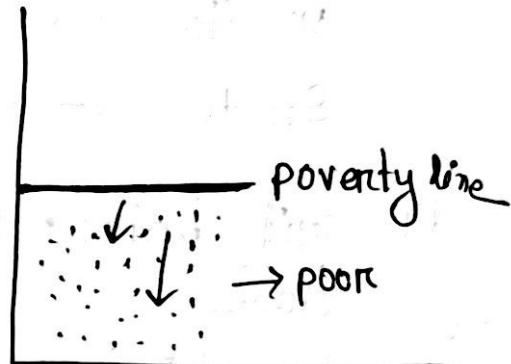
$$n = \left(\frac{Z\sigma}{E} \right)^2 \quad (\text{Mean based})$$

$$n = \frac{pqZ^2}{E^2} \quad (\text{proportion based})$$

prob. distribution $\rightarrow x/y/z \rightarrow$

over the population

Normal, poisson, Binomial



Sampling Distribution (Sampling distribution of statistics)

$$\sqrt{\text{Var}(X)} = \text{Std}(X)$$

$$\sqrt{\text{Var}(\bar{X})} = \text{Std error}(\bar{X})$$

Central limit Theorem

$$X \sim (\mu, \sigma^2)$$

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

$$\bar{X} = \hat{\mu}$$

$$S = \hat{\sigma}$$

\wedge = estimation

Standard error = $SE = \frac{S}{\sqrt{n}} = \frac{\hat{\sigma}}{\sqrt{n}}$

- $\hat{\beta} \pm se(\hat{\beta})$
-
- 95% CI means, if I draw 100 times sample we will get 100 \bar{X} 's. Among the hundred \bar{X} 's 95% \bar{X} will be very close / represent the μ .
- $N = 35$
- $n = 5$
- Bootstrap & Jackknife \rightarrow syllabus এ নাই
 ↓ ↓
 with replacement without replacement

Mid 11-12-2024

Slide 1-5

Definition মানবে না, Math মানবে

Slide-5 \rightarrow Theoretical / Descriptive

$$E = E(\hat{\theta} - \theta) = \mu$$

কথন মুচ্যে safe side মুখন p না মুখ (50%. 50%)

ওখন highest sample size পাওয়া যায়,

Standard error দেয় বস্তু,

Z table, chi-square ,	formula sheet
-------------------------	---------------

