

Statistics

Types of data

Categorical or Qualitative data

Nominal

male/female
No hierarchy

Ordinal

bad / order
avg
good

Numerical or Quantitative data

Discrete

age / rank
35, 23

Continuous

weight / height
69.5, 72.7

Population mean $\mu = \frac{\sum_{i=1}^N x_i}{N}$ [N = no. of items in population]

Sample mean $\bar{X} = \frac{\sum_{i=1}^n x_i}{n}$ [n = no. of items in sample]

- ↳ Mean is prone to outliers, median can solve outliers issue
mode can work best for categorical columns to find highest frequency
Weighted mean we assign some weights & trimmed mean we remove fraction of values from both the sides to remove outliers.

↳ Variance $\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$ used in finding dispersion of data

Mean absolute deviation $MAD = \frac{\sum |x_i - \bar{x}|}{n}$

Standard deviation $SD = \sqrt{\sigma^2}$ this holds the same unit as data

Coefficient of variation $CV = \frac{\text{Standard deviation}}{\text{mean}} \times 100\%$ It gives the idea of spread of columns & compare

↳ Graphs for univariate columns: Categorical & Numerical

- Categorical - Frequency bar charts, ^{Pie chart} relative frequency, ^{Line chart} cumulative frequency
- Numerical - Bins histograms, normal, skew, bimodal, uniform, no pattern

↳ Graphs for bivariate columns: C & C, N & N, C & N


- Categorical & C: Contingency table or cross tab

	Pclass		
Survived	1	2	3
0	21	40	11
1	12	14	50

- Numerical & N: Scatter plots

- C & N: Bar charts with aggregate fⁿ or cross tabs

Pivot tables

↳ Quantiles :- Divide a set of numeric data into equal sized group
Quartiles :- Q_1 (25 percentile), Q_2 (50 percentile), Q_3 (75 perc) 

Quartiles - Divide a set of numbers into 4 equal parts.
Quartiles: Q_1 (25 percentile), Q_2 (50 percentile), Q_3 (75 perc)

Deciles - (10th percentile) D_1, \dots, D_9

Percentile: $P_1; P_2 \dots P_{99}$

Quantile: 5 equal parts [minimum, Q_1 , Q_2 , Q_3 , max] to describe()

↳ $IQR = Q_3 - Q_1 = \text{Interquartile range [Box of boxplot]}$

$PL = \frac{P(N+1)}{100}$
 $PL =$ desired percentile value location
 $N =$ no. of obs. $P =$ percentile rank

Ex find, 75 percentile score from data

78, 82, 84, 88, 91, 93, 94, 96, 98, 99 \rightarrow sort the values

$$PL = \frac{75}{100} (10+1) = \frac{33}{4} = 8.25 \quad \therefore 95 - 98 : 95 + (98-96) \frac{1}{4} = 96.5$$

↳ Percentile of a value :-

$$\text{Percentile rank} = \frac{x + 0.5y}{n}$$

$x \rightarrow$ no. of values below given

$y \rightarrow$ no. of values equal to given

$n \rightarrow$ total no. of values

1 2 3 n
Gm 78, 82, 84, (88), 91, 93, 94, 96, 98, 99

$$\therefore \frac{3 + 0.5(1)}{10} = 0.35 = 35\%$$

Qn How to create boxplots?

6	213	241	260	281	290	314	321	350	1500
1	2	3	4	5	6	7	8	9	10

$$Q_2 = \frac{50}{100}(11) = 5 \cdot 5 = 285.5$$

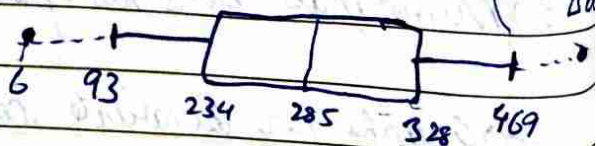
$$Q_1 = \frac{25(11)}{100} = 2.75 \therefore 213 + 0.75(241 - 213) = 234$$

$$Q_3 = 75 \times 11/100 = 8.25 = 328.25$$

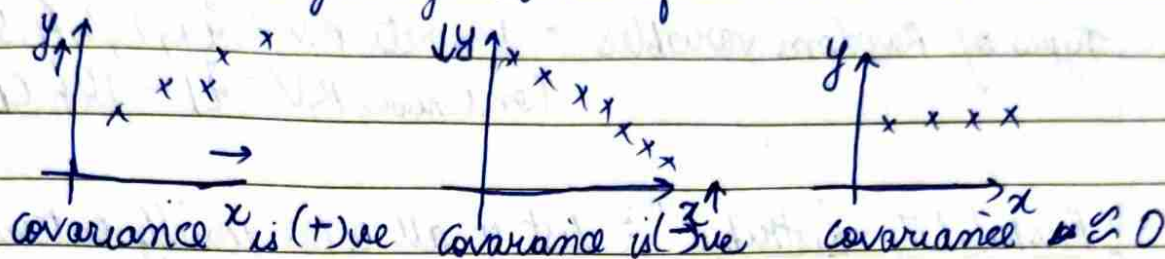
$$\text{min} = 91 - 1.5(\text{IQR}) = 93$$

$$\text{max} = Q_3 + 1.5(IQR) = 469$$

$$\text{IQR} = Q3 - Q1 = 328 - 234 = 94$$



→ Covariance : Statistical measure that describes the degree to which 2 variables are linearly related. Measures how much 2 variables change together like if one ↑ other ↑ or ↓



Covariance formula :-

Population

$$\sigma_{xy} = \frac{\sum (X - \mu_x)(Y - \mu_y)}{N}$$

$X, Y \rightarrow$ value of X, Y in population
 $\mu_x, \mu_y \rightarrow$ population mean
 $N \rightarrow$ no. of observation

Sample

$$s_{xy} = \frac{\sum (X - \bar{x})(Y - \bar{y})}{n-1}$$

$X, Y \rightarrow$ Same
 $\bar{x}, \bar{y} \rightarrow$ sample mean
 $n \rightarrow$ Same

→ Disadvantage of covariance :-

Does not tell us about strength of relatⁿ b/w 2 variable as magnitude of covariance is affected by the scale of variable.
 Covariance tells us about the spread of data. it gets scaled up

→ Correlation :- Refers to statistical relatⁿ b/w variables it measures the degree to which 2 variable are related & how they change

$$\text{Correlation} = \frac{\text{Covariance}(X, Y)}{\sigma_x * \sigma_y}$$

standard deviatⁿ of x

Indy corr

-1

0

Direct corr

1

* If 2 things are correlated does not mean it is causation or causigiet
 Causation refers to a cause & effect relatⁿ b/w 2 variables

→ Graphs for multivariate columns :-

Scatterplots 3D, Barplots with Hue, Facetgrids, Jointplots, Pairplots, Bubblecharts

→ Random variables - Set of possible values from a random exp.
 Ex coin toss H & T... $X = \{0, 1\}$ $H=0, T=1$

Random variables [capital letters]

Types of Random variables :- Discrete RV = $\{1, 2, 3, 4, 5, 6\}$ Dice
 Continuous RV = $\{1 \dots 10\}$ CGPA

→ Probability distribution - dist of all of the possible outcome of a random variable along with their probability values

Coin toss	H(0)	T(1)
Probability	1/2	1/2

We can derive mathematical fⁿ to model relatⁿ b/w outcome & prob

X → outcome 1 2 3 4 5 6

Y → probability 1/6 1/6 1/6 1/6 1/6 1/6

$y = f(x) \rightarrow$ Probability distribution function

Probability mass functⁿ (PMF)

Probability density functⁿ (PDF)

Cumulative Distribution fⁿ (CDF)

→ PMF → probability distribution of discrete random variable [dice (1-6)]

Sum of prob is equals to 1

$$y = \begin{cases} 1/6 & \text{if } x \in \{1, 2, 3, 4, 5, 6\} \\ 0 & \text{if otherwise} \end{cases}$$

CDF of PMF describes the probability of X found at $\leq x$

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

Ex for dice roll $P(4) = 1/6$

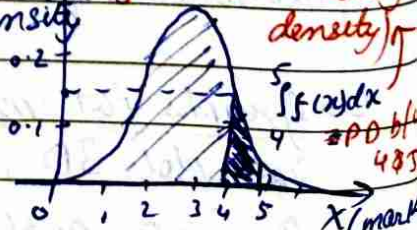
$$CMF(4) = P(X \leq 4) = P(1) + P(2) + P(3) + P(4) = 4/6$$

$$CMF(6) = P(X \leq 6) = 1 \rightarrow np \cdot \text{cumsum}(s)$$

→ PDF is a mathematical fⁿ describes prob. distributⁿ of continuous Vari

* For ∞ range of value y-axis cannot be probability Prob.

So we take small step of \int to find prob. of density getting marks b/w 4-4.001 that value is on y-axis

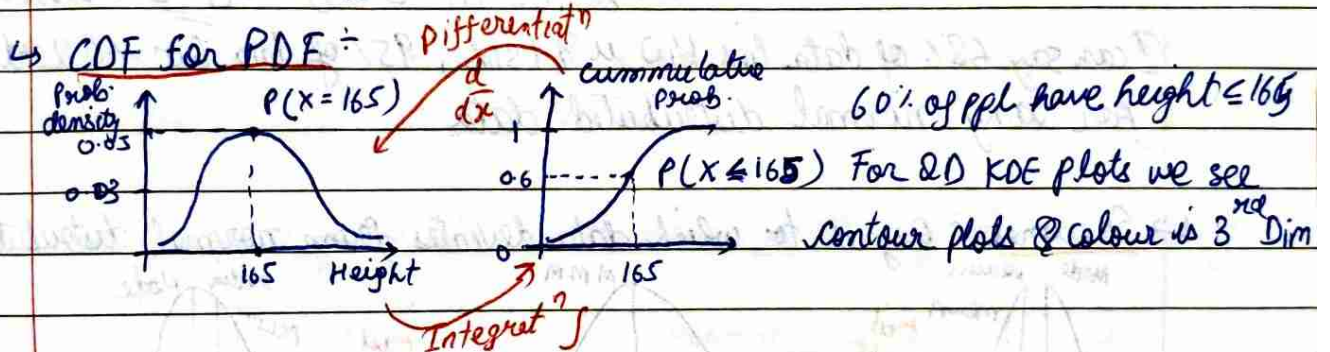
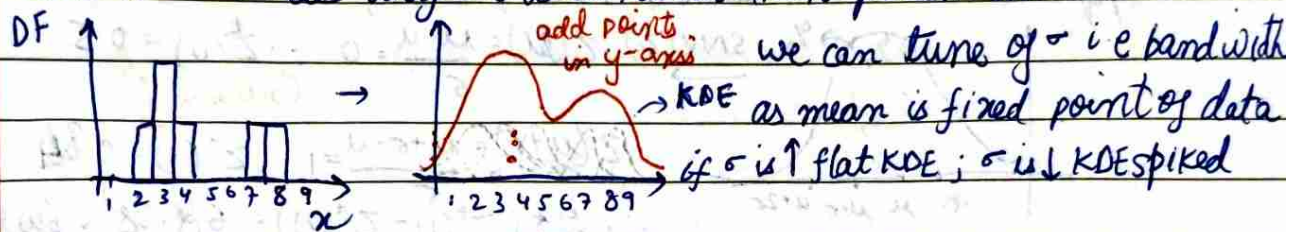


$$P(0 \leq X \leq 10) = 1 = \text{Area}$$

→ Density estimation : Statistical technique used to estimate the prob. density $f(x)$ (PDF), estimating the underlying distributⁿ of data pt.
 2 types : Parametric density estimation & Non-parametric DE
 Commonly used techniques : Kernel density estimatⁿ (KDE), histogram estimatⁿ, Gaussian mixture models (GMMs)

a) Parametric density estimation : In we have data of CRPA from 1-10 we see its histplot and observe which kind of distributⁿ say normal distributⁿ then we use the PDF formula of normal with parameters of μ & σ to create the distribution

b) Non-parametric DE : without making underlying assumptions here we will be using the data rather than parameters



→ Normal distribution : continuous probability distributⁿ symmetric around mean and in bell curve (Gaussian)

2 parameters : mean (μ) centre of distributⁿ PD

Standard deviatⁿ (σ) spread of distributⁿ

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

belongs to Normal

PDF Eqⁿ of normal dist asymptotic

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

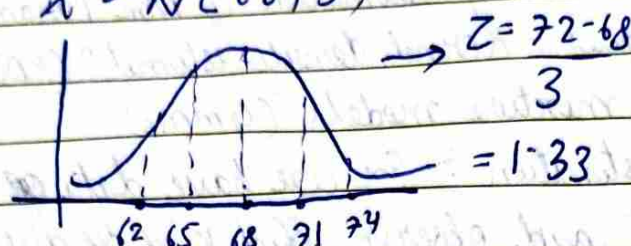
Standard normal variate $Z \sim N(0, 1)$ mean $\mu = 0$ & std. $\sigma = 1$

Convert any column to standard normal variate

$N(\mu, \sigma)$ Age $\rightarrow Z \sim N(0, 1)$ Symmetrical
 $z_1 = \frac{27 - \mu}{\sigma}$ mean = median = mode
 $z_2 = \frac{61 - \mu}{\sigma}$ 68-95-99.7 Rule
 Age is 1

- Q) Height of adults follow normal distribution with $\mu = 68$ & $\sigma = 3$ inches. Find probability that a random male is taller than 72 inches?

$$\therefore X \sim N(68, 3)$$



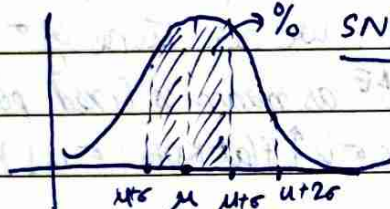
from Z-table

$$Z(1.33) = 0.9098$$

$$\therefore 1 - 0.9098 = 0.0902$$

Z-table of std. normal dist. is computed \therefore 9% chance tall > 72 inches

- Q) For a normal distribution $X \sim N(\mu, \sigma)$ what percent of population lie b/w mean & 1std, 2std, 3std?



$$Z(\mu) = \frac{\mu - \mu}{\sigma} = 0 \therefore Z(0) = 0.5$$

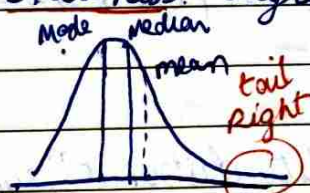
Area $\int_{-\infty}^x$

$$Z(\mu + \sigma) = \frac{\mu + \sigma - \mu}{\sigma} = 1 \therefore Z(1) = 0.84$$

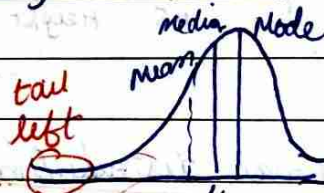
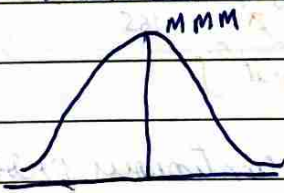
$$\therefore 2 \times (Z(1) - Z(0)) = 68.28\% \text{ b/w } \mu \text{ & } 1\sigma$$

I can say 68% of data lies b/w μ & 1std, 95% of data b/w μ & 2std for any normal distributed data.

\hookrightarrow Skewness: Degree to which data deviates from normal distribution



Right skew / (+) skew



Left skew / (-) skew

Greater the skew distance b/w mean, median, mode

* Statistical moment: 1st moment = mean

2nd moment = Variance, 3rd moment = Skew, 4th moment = Kurtosis

$$\text{skewness} = \frac{n}{(n-1)(n-2)} \sum \left(\frac{x - \bar{x}}{s} \right)^3$$

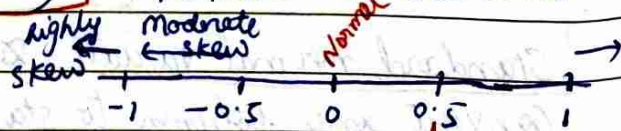
Sample skew

pandora uses this skew

\hookrightarrow CDF on Normal curve

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t-\mu}{2\sigma^2}} dt$$



No skew \approx Normal

skew ≈ 0

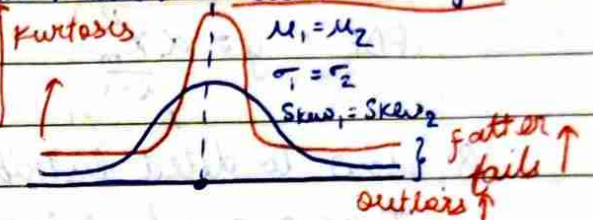
not always for normal graph

↳ Used in data science to detect outliers, hypothesis testing, Central limit theorem, assumptⁿ of ML algo like GMM & linear Reg (errors is normally distributed)

↳ Non-gaussian distribution → Continuous non-gaussian
↳ Discrete non-gaussian

Kurtosis is 4th statistical moment, it measures tailedness of PD

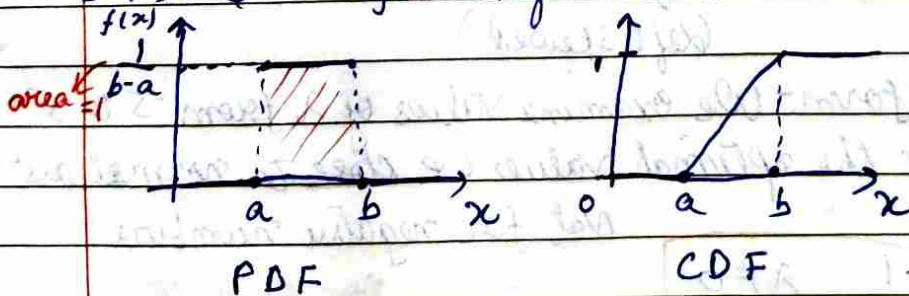
$$Kurtosis = \frac{n \times (n+1) \times \sum_{i=1}^n (x_i - \bar{x})^4}{(n-1)(n-2)(n-3) \times s^4} - \frac{3 \times (n-1)^2}{(n-2)(n-3)}$$



↳ QQ plot : (quantile-quantile) plots checks for distribution of data how ^{similar} normal it is. Take a column X and creates its 1st quantile, 2nd ... 99th quantile and do the same for a normal data. Now draw a scatterplot and if it fits in same line then data is normal and gaussian in nature. It can use any other distribution and can compare how similar it is to our data.

↳ Uniform distribution : All outcomes are equal within range
Uniform → Continuous → dice {1, 2, 3, 4, 5, 6}
Uniform → Discrete
Denoted as $X \sim U(a, b)$ → parameters $a = 1$ lower, $b = 6$ higher

↳ PDF & CDF for uniform distribution (skewness = 0)



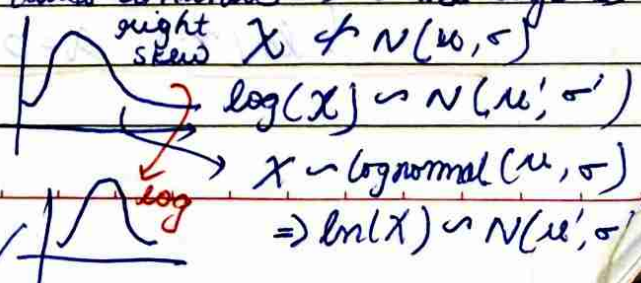
PDF of lognormal

$$\frac{1}{x \sigma \sqrt{2\pi}} e^{-\frac{(\ln(x) - \mu)^2}{2\sigma^2}}$$

↳ Log normal distributⁿ : Heavy tailed continuous PD whose logs is normally distributed.

* Distributⁿ jiska log is normal is called as log-normal

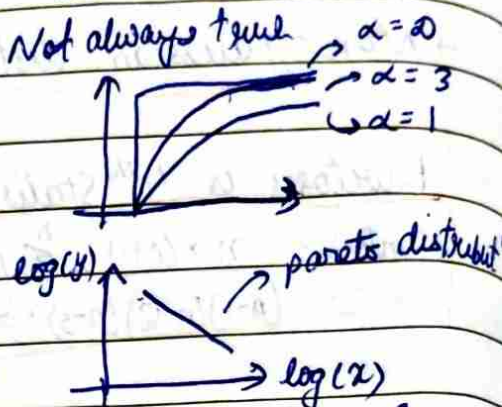
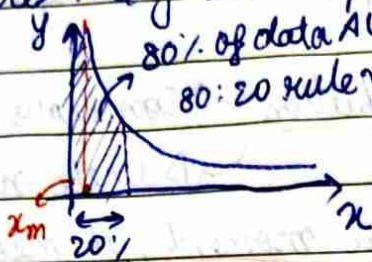
Say if we have right skew data apply log & the value are normal. ✓



↳ Pareto distribution : Exhibits power-law behaviour
 power law where x & y are exponentially proportional

$$y = k \cdot x^{\alpha}$$

$\alpha \uparrow$ peak \uparrow tail \downarrow
 $\alpha \downarrow$ peak \downarrow tail \uparrow



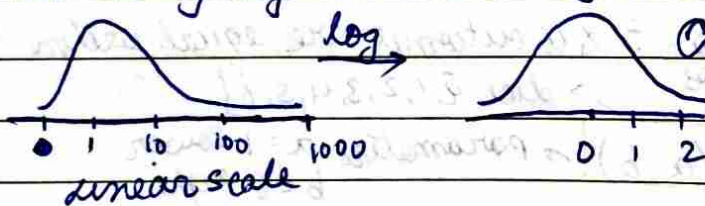
$$\text{PDF} = y = \frac{\alpha x_m^{\alpha}}{x^{\alpha+1}}$$

Q) How To detect distributⁿ is pareto?

Take a graph of $\log(y)$ vs $\log(x)$ is straight line ✓
 We can also use QQ plot take our reference as pareto plot

↳ Transformation : To convert data into normal distributⁿ
 f^{trans} → log transform, reciprocal transform, power (Sq, Sqrt),
 power tr → Box-cox transform, Yeo-johnson transform.

a) log transform : Take log of values, if data is right skew
 On taking log we break the interval into smaller ones.



b) Reciprocal ($1/x$) (c) Square (x^2) (d) Sqrt (\sqrt{x})
 (left skewed)

e) Box-cox transform : We examine values of λ from -5 to 5
 and we choose the optimal value i.e. close to normal dist
 Not for negative numbers.

$$x_i^{\lambda} = \begin{cases} \frac{x_i^{\lambda} - 1}{\lambda} & \lambda \neq 0 \\ \ln(x_i) & \lambda = 0 \end{cases}$$

f) Yeo-Johnson transform - Adjusts box-cox & can apply on (-)ve numbs

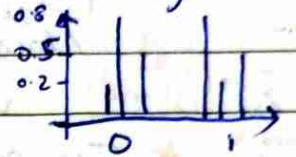
$$x_i^{\lambda} = \begin{cases} [(x_i+1)^{\lambda} - 1] / \lambda & \text{if } \lambda \neq 0, x_i \geq 0 \\ \ln(x_i+1) & \text{if } \lambda = 0, x_i \geq 0 \\ -[(1-x_i+1)^{2-\lambda} - 1] / (2-\lambda) & \text{if } \lambda \neq 2, x_i < 0 \\ -\ln(-x_i+1) & \text{if } \lambda = 2, x_i < 0 \end{cases}$$

↳ Bernoulli Distribution - Binary outcome, either success or failure

PMF = $P(X=x) = p^x (1-p)^{1-x}$; p = prob of success

$P(X=0) = 0.2$ & $P(X=1) = 0.8$ $1-p$ = prob of failure

$P(X=0) = 0.5$ & $P(X=1) = 0.5$



for binary classificⁿ like spam/ham, fraud etc falls in Bernoulli

↳ Binomial distribution - ^{tossing coin 3 times (independent events)} Number of successes in a fixed number of ^{independent} Bernoulli trials, with 2 possible outcomes (success & failure) where probability of success is constant for each trial. It has 2 parameters n → no. of trials & p → probability of success.

8) Probability of watching lecture is 0.5 (3 ppl) → $\frac{2}{2} \frac{2}{2} \frac{2}{2} = 2^3 = 8$ outcomes
No-one watches out of 3 people = $1/8$ (NNN)

1 out of 3 watch = $3/8$ [YNN, NNY, NYN] But what if we can't get

2 out of 3 watch = $3/8$ [YYN, NYY, YNY] the sample space

3 out of 3 watch = $1/8$ [YYY] we will use the formula

PDF of Binomial = $P(X=x) = {}^n C_x p^x (1-p)^{n-x}$ n = no. of trials

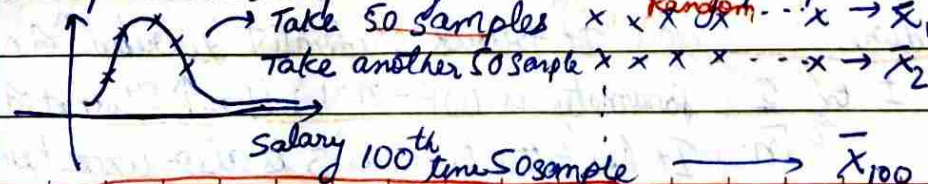
p = prob. of success

x = desired result

→ ${}^3 C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^1 = \frac{3}{8}$

↳ Sampling distribution - It is PD that describes the statistical property of a sample statistic (such as mean, variance) for multiple

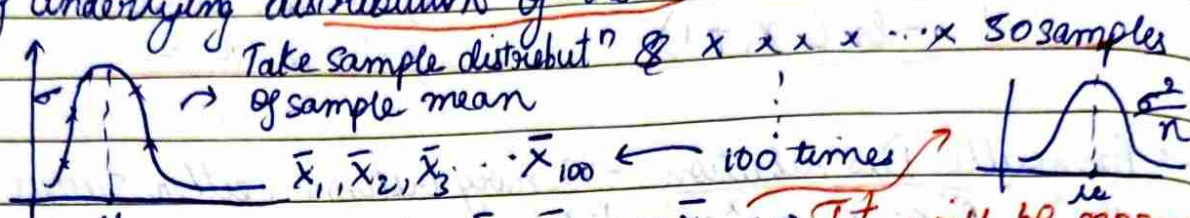
independent samples of same size or proportion. ^{Estimate the variability of a sample statistic}



→ $\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_{100}$ → Sampling distribution of sample mean

Take mean/Variance

- ★ Central limit theorem - CLT states distribution of sample means of large no. of independent & identical distributed random variable will approach to normal distribution, regardless of underlying distribution of the variables.



Now if we plot $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_{100} \rightarrow$ It will be normal irrespective of original distrib^n, skewed, lognormal, etc.

\bar{x}_1 point estimate
single value

- ★ The mean will be the same, the variance will reduce by a factor of n i.e. no. of samples. σ^2/n [sample $n > 30$]
- ★ Now we cannot say that μ of sample mean is same or approx to population mean. We can say about our confidence interval. and say we are 95% sure that our mean will lie between $[\mu - 2\sigma, \mu + 2\sigma]$ as this holds 95% of our value. so we are sure about mean.

- ★ Ways to calculate CI - \rightarrow Z-procedure [pop \rightarrow std available σ]
 \rightarrow t-procedure [pop \rightarrow std X available not]
Confidence intervals is created for parameters & not statistics. Statistics help us get the confidence interval for a parameter.

- ★ Find confidence interval using Z-procedure:

Condition \rightarrow Z-procedure formula used where data is randomly selected, standard deviation is known for a population, data is normal.

$$CI = \text{point estimate } (\bar{x}) \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$1 - \alpha = \text{confidence level}$

$1 - \alpha = 95\%$

$\sigma = \text{std population}$

$n = \text{sample size} \Rightarrow \text{say } 100$

table at 95% confidence (Critical Value)

table at 95% generally

- ★ t-procedure - It follows student's T distribution \approx Normal. Parameter is $DOF = n - 1$; if $n \uparrow$ Student \approx Normal. It has fatter tails as s is also uncertain so range increases of CI; $t_{\alpha/2} > Z_{\alpha/2}$ But $n \uparrow t_{\alpha/2} \approx Z_{\alpha/2}$

$$CI = \bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$$

point est. sample mean

no. of sample

↳ Hypothesis testing - A method of statistical inference used to decide whether the data at hand sufficiently support a particular hypothesis. It allows us to make probabilistic statements about population parameter.

Ex: Lays has green packet chips if they substitute it with say purple so will there sales ↑, what is the prob of this hypothesis

H_0 ↳ Null hypothesis (Status quo) - It is a statement which assumes no significant change. It serves as starting point. It is the assumption of no effect until proven otherwise i.e. lawyer proves that he is the guilty one. & null hypothesis (H_0) say he is not guilty / not committed crime. The purpose of hypothesis testing is to gather evidence to reject or fail null hypothesis and prove in the favour of alternative hypothesis (H_1 or H_a)

(H_1 or H_a) ↳ Alternative hypothesis (Research hypothesis H_1 or H_a) - It contradicts null hypothesis and claims significant result. We need to contradict and prove the man is accused. Always any one of them is true either null or alternative hypothesis. We try to gather evidence to reject the null hypothesis.
⇒ Failing to reject null hypothesis does mean null hypothesis is true. We couldn't gather much evidence to prove him accused. That does not mean he is completely not guilty.

↳ Rejection region approach:

1. Formulate a Null and alternate hypothesis
2. Select significant lvl (^{95%} 0.05 or ^{99%} 0.01) of rejecting null hypothesis
3. Check assumptions (example distribution) → Normal, Lognormal, categories
4. Decide appropriate test (Z-test, T-test, Chi-square test, ANOVA)
5. State the relevant test statistic ^{If (n) pop is known} → for categorical columns
6. Conduct the test
7. Reject or not reject null hypothesis.
8. Interpret the result.

- 8) Suppose company has training program. Avg. productivity before program is 50 units with a std. of 5 units. After programme we took sample of 30 employees and productivity increased by 53 units. How significantly productivity has increased?

1) $H_0 = \mu = 50$ Null hypothesis $\bar{x} = 53$ $\sigma = 50$
 2) $H_a = \mu > 50$ Alternate hypothesis $n = 30$

3) $\alpha = 0.05 \rightarrow 5\%$

4) Normality is valid as sample ≥ 30 / pop std. (σ) is known

5) Z-test

6) $Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{53 - 50}{5 / \sqrt{30}} = 3.28$

Std normal variate σ / \sqrt{n}

No reject region

Reject region

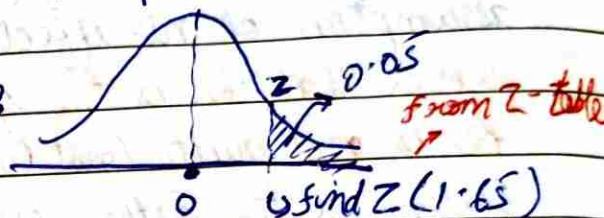
0

1.65

3.28

(7) reject H_0

AS The Z has fallen in reject region so we can reject the null hypothesis with 95% confidence



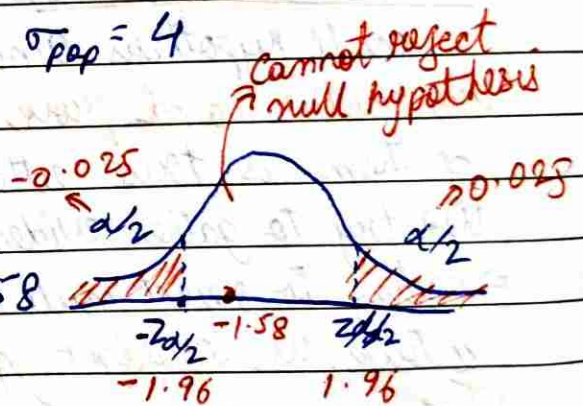
9) $\mu = 50$, $n = 40$, $\bar{x} = 49$, $\sigma_{pop} = 4$

1) $H_0 = \mu = 50$, $H_a = \mu \neq 50$

2) $\alpha = 0.05$

3) Normality \checkmark σ \checkmark Ztest \checkmark

4) Z test, $Z = \frac{49 - 50}{4 / \sqrt{40}} = -1.58$



1) Type I v/s Type II Error: Hypothesis $H_0 \times \rightarrow \checkmark$

Type I is false positive; rejects H_0 when \checkmark $H_0 \checkmark \rightarrow \times$

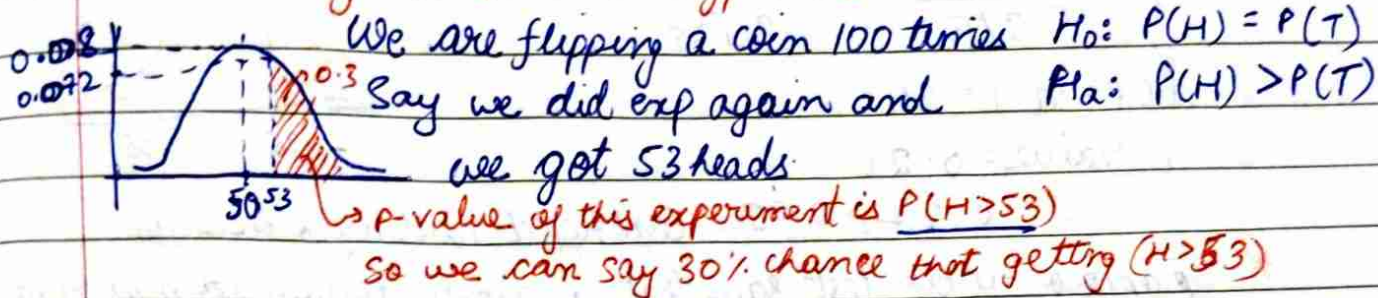
(2) Type II is false negative; accepts H_0 when \times

Prob. of committing Type II error is β

Decision	H_0 true	H_a true
Reject H_0	Type I	\checkmark
Accept H_0	\checkmark	Type II

Hypothesis testing is used to compare ML models using CV-folds. Used in feature selection using t-test, χ^2 , ANOVA, hyperparameter

↳ p-value - Probability of getting a sample as more extreme (having more evidence against H_0) than our own sample given the null hypothesis is true. It measures the strength of evidence against the Null Hypothesis.



★ If we have $P\text{-value} \leq \alpha$ we reject your null hypothesis here α is generally 0.05 so $0.3 \neq 0.05$ so we don't have enough evidence to reject our null hypothesis. If $P(H > 80)$ we get p-value ~ 0 so we can reject the H_0 and say coin is fair.

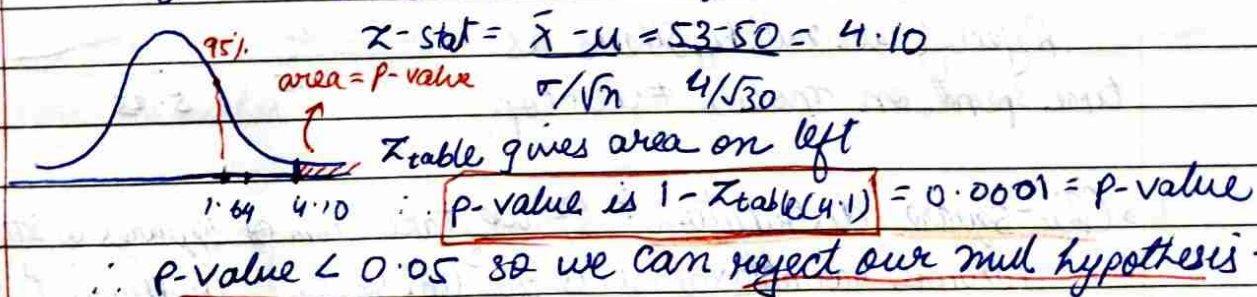
↳ Without significance value (α):

Very small p-value ($P < 0.01$) strong evidence against H_0
 Small p-value ($0.01 < P < 0.05$) moderate evidence against H_0
 Large p-value ($0.05 \leq P < 0.1$) weak evidence against H_0
 Very large p-value ($P \geq 0.1$) weak or no evidence against H_0

8) P-Value in context of Z-test

$\mu = 50$, $n = 30$, $\bar{x} = 53$, $\sigma_{pop} = 4$; find if there's significant \uparrow ?

$\therefore H_0: \mu = 50$ & $H_a: \mu > 50$, $\alpha = 0.05$ (one tailed test)



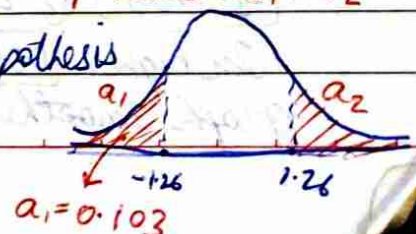
8) $\mu = 50$, $n = 40$, $\bar{x} = 49$, $\sigma = 5$, $\alpha = 0.05$ | $H_0: \mu = 50$ & $H_a: \mu \neq 50$ (2-tailed test)

$Z = 49 - 50 = -1.26$

$5/\sqrt{40}$

\therefore p-value = 0.206

Can't reject our null hypothesis



Q) $\mu = 50$, $n = 25$, $\bar{X} = 49.7$, $S = 1.2$, $\alpha = 0.05$ pop σ is X

$H_0: \mu = 50$, $H_a: \mu \neq 50$

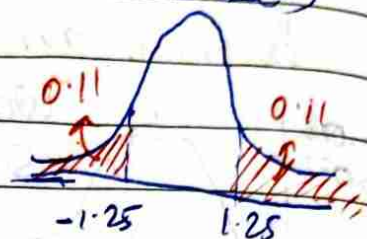
Assuming sample is normal; we use t-test here as (3)

$$t = \frac{\bar{X} - \mu}{S/\sqrt{n}} = \frac{49.7 - 50}{1.2/\sqrt{25}} = -1.25$$

$\text{dof} = n - 1 = 24$

t-value = 0.22

$\therefore 0.22 > 0.05 \therefore$ Can't reject null hypothesis H_0



* Shapiro-Wilk test says that a sample follows normal or not from scipy. stats import shapiro \rightarrow returns statistic & p-value if p-value < 0.05 not normal | p-value > 0.05 normal

\hookrightarrow Independent 2 sample t-test

like we can test if males avg. age $>$ or $<$ females avg. age

\hookrightarrow Independent obs., normally distributed, equal variance homoscedasticity
if $\sigma_A^2 \neq \sigma_B^2$ checked by F-test, if $\sigma_A^2 \neq \sigma_B^2$ use Welch's t-test ($\sigma_A^2 = \sigma_B^2$)

\hookrightarrow here eq. not required

Q) $H_0: \mu_d = \mu_m$

$H_a: \mu_d \neq \mu_m$

$\text{dof} = n_1 + n_2 - 2 = 58$

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

$n_1 = 30$

$n_2 = 30$

$\bar{X}_1 = 18.5$

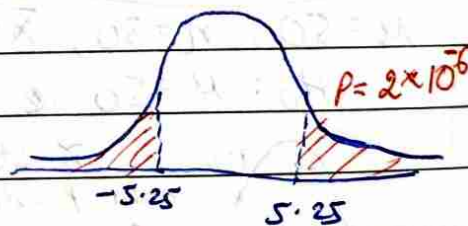
$\bar{X}_2 = 14.3$

$\sigma_1 = 3.5$

$\sigma_2 = 2.7$

$$t = \frac{18.5 - 14.3}{\sqrt{(3.5)^2/30 + (2.7)^2/30}} = 5.25$$

Rejecting our null hypothesis as time spent on mobile \neq desktop



\hookrightarrow Chi-square distribution: If we take sum of squares of std. random normal variable values is χ^2 (chi square), continuous P.D.

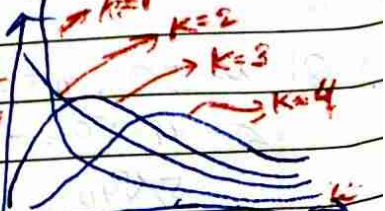
$X_1 \sim N(0,1) \therefore Q_1 = Z_1^2 = X_1^2$

$Z_2 \sim N(0,1) \therefore Q_2 = Z_1^2 + Z_2^2 = X_2^2$

$Z_3 \sim N(0,1) \therefore Q_3 = Z_1^2 + Z_2^2 + Z_3^2 = X_3^2$

On higher $\text{DOF} = k$ the graph smoothens $\hookrightarrow k$

Sum of squares of a normal value will be 0 if $k=1$ as major values lie b/w -1 to 1 where square is close to 0.



↳ Chi-square test : We use this test to check for significant ^{non-parametric} association between categorical variables. We can check goodness of fit test as well if observed distribution of a single categorical variable matches an expected theoretical distribution.

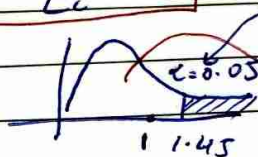
Q) Dice is rolled 60 times and we need to check whether dice is fair

H_0 : dice is fair \rightarrow uniform distribution $\rightarrow 10 \mid 10 \mid 10 \mid 10 \mid 10 \mid 10$

H_1 : dice is unfair \rightarrow Observed distribution $\rightarrow 12 \mid 8 \mid 11 \mid 9 \mid 10 \mid 10$

Formula $\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$ = Here $K=6$ as 6 outcomes
 $\alpha = 0.05$ $\therefore D.O.F = K-1 = 5$

$$\chi^2 = \frac{4+4+1+1}{10} = 1$$



We will check χ^2 values for $D.O.F=5$ from chi-table

\therefore we cannot reject H_0

Dice is fair

↳ F-distribution : Continuous PD used in statistical hypothesis and analysis of variance (ANOVA). It has 2 parameters df_1 & df_2 . Its left bound is at zero, (+)ve skewed. Used to compare fits of different statistical models. F-statistic is calculated by dividing the ratio of two sample variance or mean squares from ANOVA table. Value is compared from F-distribution for its confidence.

↳ One way ANOVA test :

H_0 : all the group means is equal

H_1 : Atleast one group mean is diff

Now its just a creation of big table which computes b/w sample or categories & within sample or categories

$$F = \frac{\chi_1^2 / d_1}{\chi_2^2 / d_2}$$

