

StatsBasic

① Variable

② Rand "

③ Population, Sample

" mean, " mean

④ Population distribution

Sample

Sampling

⑤ Mean, median, Mode

⑥ Range

⑦ Measure of Dispersion

⑧ Variance

⑨ STD

⑩ Gaussian / Normal distribution,

Intermediate

⑪ Standard Normal distribution

⑫ Z score

⑬ Prob density function

⑭ Cumulative distribution function

⑮ Hypo testing

⑯ Graphs

⑰ Kernel-Density estimation

⑱ Central limit theorem

⑲ Skewness of data

⑳ Covariance

㉑ Pearson corr coeff

㉒ Spearman Rank corr

㉓ Hypo test

Advanced

㉔ Q-Q Plot

㉕ Chebyshov's inequality

㉖ Discrete and cont distib.

㉗ Bernoulli and Binomial

㉘ Log Normal Distribution

㉙ Pois & w distribution

㉚ Box Cox transform

㉛ Poisson Distribution

㉜ App of non fause Distr

Start

Descriptive Stats

Inferrential - Stats

* Inf - Stats: is a statistical method that deduce from a small but representative sample the characteristic of a bigger population.

Ex. Asking people of diff election areas and conclude which party might win. By statistical testing to be conducted.

Ex of inferential stats that we do

→ hypothesis testing

① Z-test

② T-test

③ Chi Square test

$$\text{mean } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

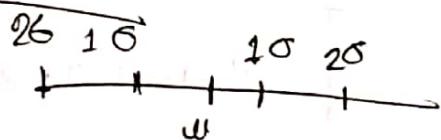
$$\text{Var} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\sigma = \sqrt{\text{Var}}$$

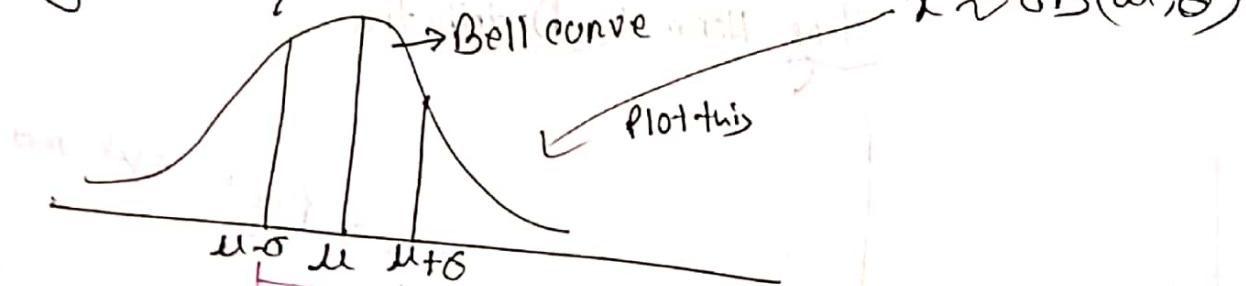
Variance and std help to understand, from the mean, how far the values are distributed.

Gaussian / Normal distribution

Consider x as random variable



x belongs to Gaussian distribution with some mean and sigma value / standard deviation.



Empirical formulas: Probability that 68% of data points belonging to the random variable x fall within the range of first STD.

$$① P_n[\mu - \sigma \leq x \leq \mu + \sigma] \approx 68\%$$

As H0 is Normal distribution

$$② P_n[\mu - 2\sigma \leq x \leq \mu + 2\sigma] \approx 96\%$$

$$③ P_n[\mu - 3\sigma \leq x \leq \mu + 3\sigma] \approx 99.7\%$$

Use of them to remove outliers

Using std:

$$\text{std} = \text{df.column.std} = 5$$

$$\text{df}[(\text{df.column} \leq 3 * \text{std}) | (\text{df.column} \geq 3 * -\text{std})]$$

Using z score:

$$z\text{ score} = \frac{x - \bar{x}}{\sigma}$$

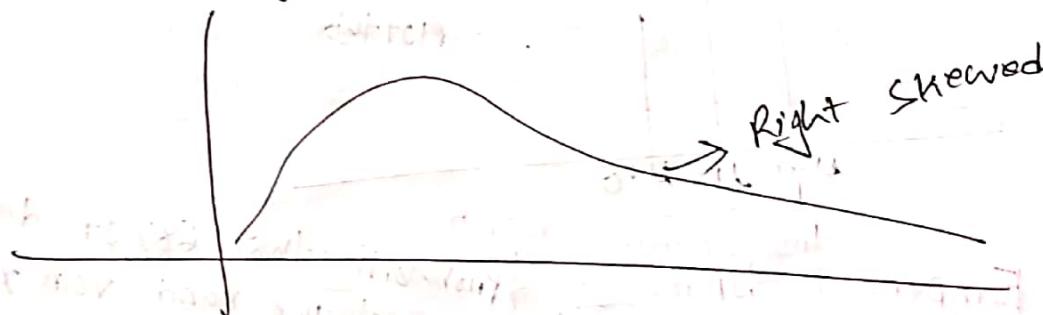
* Remove if z score is greater than 3.
It is similar to std.

To remove outliers from dataset

$$df_{no_outlier} = df[(df.zscore < 3) \& (df.zscore > -3)]$$

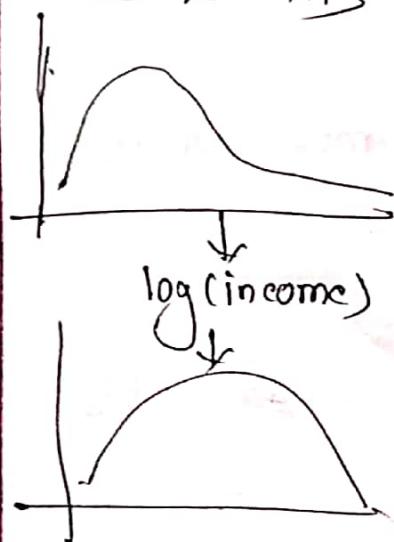
We can also use percentile method to remove outliers.

Log Normal Distribution



This is used to normalize the data.

for the above graph, if we apply $\log(x)$ to the x axis, it becomes normal distribution.



If you get a normal distribution by applying a log function to a dataset the dataset is log normally distributed.

$$\text{Gaussian} \rightarrow \text{Standard Normal distribution}$$

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

Normalization $x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$

\Rightarrow Standardization

$$\text{std, } \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$

(3)

* Normalization or Standardization.

\Rightarrow Normalization, when data doesn't follow Gaussian distribution $\frac{x - x_{\min}}{x_{\max} - x_{\min}}$ mean & std are not constant

\Rightarrow Standardization, when data follows Gaussian distribution. $\frac{x - \bar{x}}{\sigma}$ \Rightarrow It converts the data as mean $\rightarrow 0$ and std $\rightarrow 1$

* It maps relationship between two features.

$$\text{Cov}(\text{size}, \text{price}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$\downarrow \quad \downarrow$

It is similar to variance

$$\begin{aligned} \text{Var}(x) &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x}) \end{aligned}$$

$$\text{So, } \text{Cov}(x, x) = \text{Var}(x)$$

It just says that there is positive or negative relationship but doesn't say about how much positive / negative.

important
for feature
selection

Pearson Correlation Coefficient

① Strength

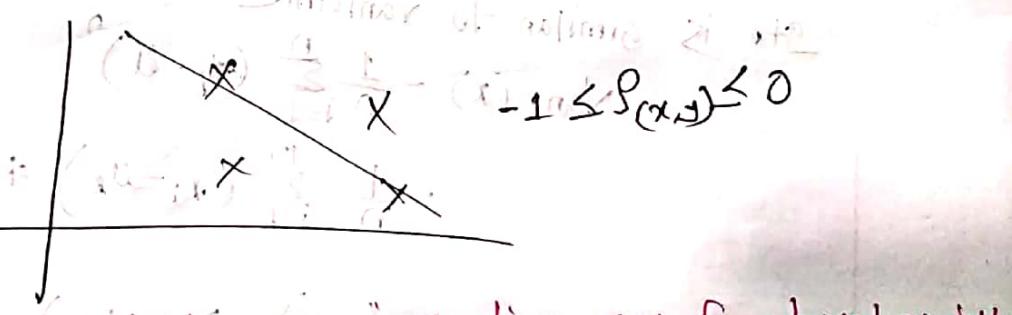
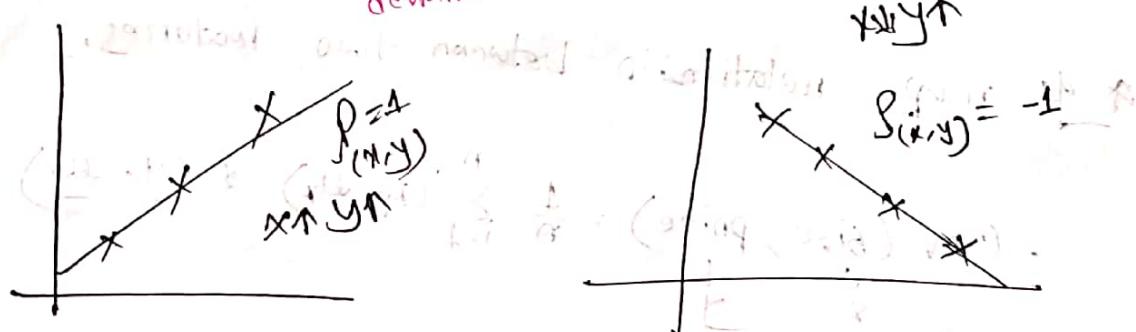
② Direction of Relationship.

$$\text{Range} = [-1 \text{ to } +1]$$

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

$$\text{formula } \rho_{(x,y)} = \frac{\text{cov}(x,y)}{\sigma_x * \sigma_y} \rightarrow \text{covariance}$$

$$\text{Standard deviation} \rightarrow \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}}$$



① find out which variables are impacting most for target variable

② find out multicollinearity and remove them.

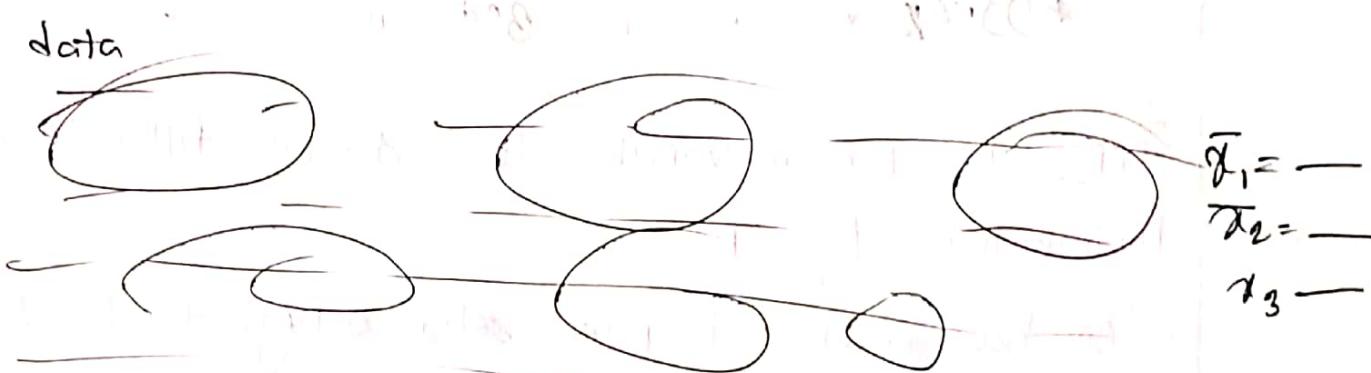
without target variable → different features are highly correlated

The presence similar info, we can remove them

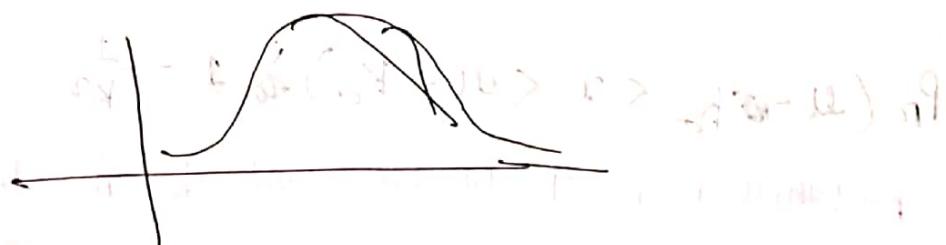
(4)

Central Limit theorem

* No matter population distribution, if we take enough sample, calculate the mean of them, the distribution of the mean will be normal/gaussian.



If we plot $\bar{x}_1, \bar{x}_2, \bar{x}_3$, distribution would look like this:



And their mean will be closer to population mean. And standard deviation will be n times less

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

$$\text{standard deviation} = \sqrt{\frac{1}{n}} \cdot \sigma = \frac{\sigma}{\sqrt{n}}$$

CHEBYSHEV'S INEQUALITY

In Gaussian distribution;

Suppose y is random variable that doesn't follow Gaussian Distribution $y \not\sim \text{GD}$.

So, how much data points ~~between~~ belongs to 1st STD

$$\Pr(\mu - \sigma < x < \mu + \sigma) \geq 1 - \frac{1}{K^2}$$

$$\Pr(\mu - \sigma k_0 < x < \mu + k_0) \geq 1 - \frac{1}{k^2}$$

k = which range of std, we want to find the value.

Suppose, $k=2 \rightarrow$ for 25D

$$P_n(\mu - 2\sigma < \bar{X} < \mu + 2\sigma) \geq 1 - \frac{1}{2^n}$$

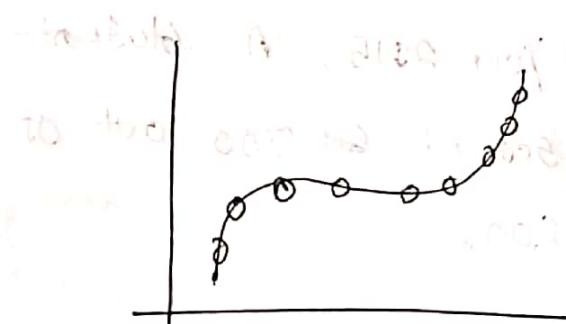
for $k=3$ for 35TD

$$P_n > 1 - \frac{1}{3^2} = 1 - \frac{1}{9} = \frac{8}{9} \approx 85\%$$

(5)

better than Pearson correlation
As more accurate on
NON LINEAR DATA

Spearman's rank correlation Coefficient



Pearson correlation

$$= \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

Spearman correlation = 1

Pearson correlation $r = 0.88$ (not highly correlated)

formula, $R_s = \frac{\text{cov}(R_{gx}, R_{gy})}{\sigma_{R_{gx}} \cdot \sigma_{R_{gy}}} = \frac{\text{cov}(R_{gx}, R_{gy})}{\sigma_{R_{gx}} \cdot \sigma_{R_{gy}}}$

$R_{gx} = \text{Rank}_x - \bar{\text{Rank}}_x$

x	y
5	7
1	2
4	3
2	1

Sort

x _{rank}	y _{rank}	Rank _x	Rank _y	d _i	d _i ²
1	2	1	2	-1	1
2	1	2	1	1	1
4	3	3	3	0	0
5	7	4	9	0	0

If all n ranks are distinct integers, then

Spearman's, $R_s = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$

$$= 1 - \frac{6 \times 2}{9(16-1)} = 1 - \frac{12}{60} = \frac{60-12}{60}$$

$$= \frac{48}{60}$$

$$= 0.8 \rightarrow \text{highly correlated.}$$

Quantile & Percentile

Year 2011, A student scored 450 out of 600

Year 2015, A student scored ~~450~~ 500 out of 600.

Who is better? Can't decide as year is different.
The question might be higher.

Better comparatively

We need more info

82% students got less than 450.

73% students got less than 500.

Ex. The 63rd percentile of the mark in mathematics in a class is 82.

This means 63% students of that class got 82 or less

⇒ Higher the percentile, rarer your data.

* There are total 99 percentile.

90th percentile is also called upper 10th percentile.

(6)

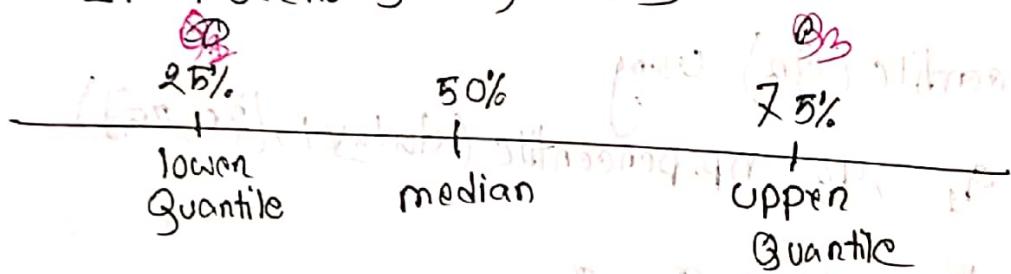
Quantile

Percentile is a version of quantile

Percentile considers 100 equal sections to locate our data

Quantile can consider any number of equal sections

If 4 sections = Quantiles



Calculated Interquartile Range

$$Q_1 \text{ to } Q_3$$

$$IQR = Q_3 - Q_1$$

Data 3, 5, 7, 1, 1, 8, 9, 6

↓ sorted

$$1, 1, 3, 5, 6, 7, 8, 9$$

$$Q_1 = 2 \quad Q_3 = 6.5$$

$$\text{Interquartile Range} = Q_3 - Q_1 = 6.5 - 2 = 4.5$$

Outlier Removal

Two methods:

① Z-score

② Interquartile Range ($75\% - 25\%$)

Interquartile Range

We can find lower percentile (Q_1) and upper percentile (Q_3) using

$q_1, q_3 = np.percentile(\text{dataset}, [25, 75])$

$$\text{iqr} = q_3 - q_1$$

$$\text{lower-bound-value} = \text{quantile} - (1.5 * \text{iqr})$$

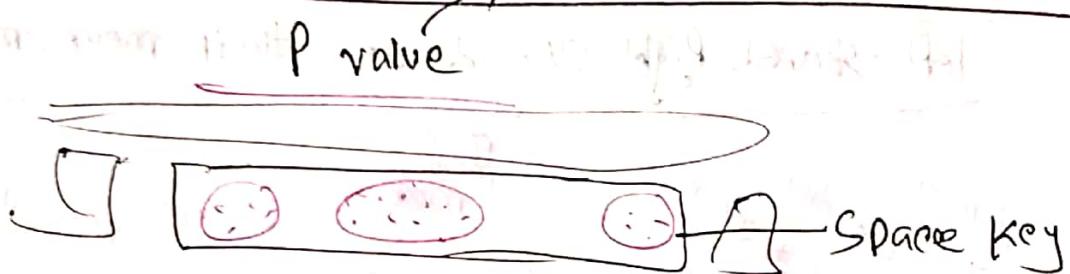
$$\text{upper-bound-value} = \text{quantile} + (1.5 * \text{iqr})$$

We can filter the values between lower and upper bound value and remove rest to remove outliers.

Z-score: provides insight about how much standard deviation away from the mean. We can remove data where Z score is -3 and $+3$ less than -3 and more than $+3$.

(4)

Used to accept/reject null hypothesis



Experiment on type on keyboard.

If we repeat this experiment 100 times and find that stroke on this place one time



P-value: Probability that the null hypothesis is true

Ex: Toss a coin 100 times \rightarrow 50 times it is head.

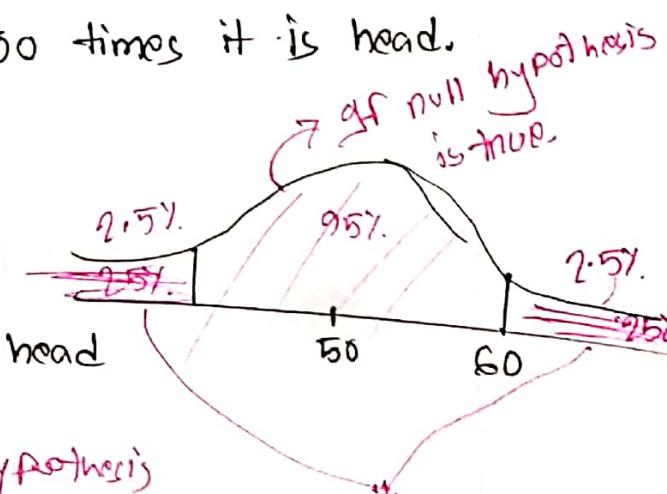
H_0 = The coin is fair.

H_1 = The coin is not fair.

\Rightarrow Suppose we get 60 times head

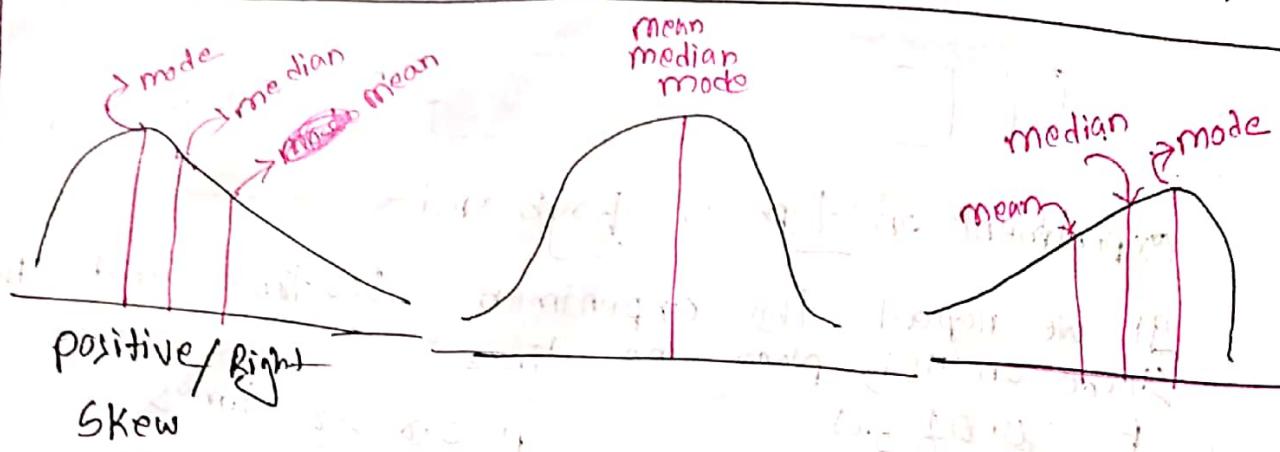
If $P \leq 0.05$, reject null hypothesis

as it is away from the mean value.



If they fall in this region we'll reject hypothesis and p value will be less than 0.05

left-skewed, right skewed, and their mean, median, mode



right-skewed: have only one offidately higher

most of the numbers are small and mode is large

most of the numbers are off

↳ normal distribution, Standardized normal distribution
↳ normal distribution

(8)

Bernoulli distribution

Two outcome: ex: Tossing a coin.

$X=1 \rightarrow$ Head
 $X=1 \rightarrow$ Success
 $X=0 \rightarrow$ Tail
 $X=0 \rightarrow$ Fail

Prob of (X) will have some value x

$$P(X=x) = P^x (1-P)^{1-x} \rightarrow \{ \text{probability mass function} \}$$

Possible value for X is

$$X=0 \quad \text{or} \quad X=1$$

Two examples of bernoulli distribution

$$P(X=0) = 0.2 \text{ and } P(X=1) = 0.8$$

$$P(X=0) = 0.8 \text{ and } P(X=1) = 0.2$$

$$2.0 = 1.0 - 1.0 = (-1)^0$$

$$P(\text{success}) = P(\text{head})$$

$$P(\text{failure}) = 1 - P = q$$

We know $P(X=x) = P^x (1-P)^{1-x}$
failure:

$$\Rightarrow P(X=0) = P^0 (1-P)^{1-0}$$

$$= P^1 - P$$

$$= 0.8$$

$$= 0.2$$

success: 0 :

$$P(X=1) = P^1 (1-P)^{1-1}$$

$$= \underline{\underline{P}}$$

$$\text{PMF} = \begin{cases} \sigma = 1 - p & \text{if } x=0 \\ p & \text{if } x=1 \end{cases}$$

Supported by $P^x (1-p)^{x-1} \Rightarrow \text{PMF}$

Probability for fair coin:

$$P(H) = p = 1/2$$

~~$P(T) = 1 - P(H) = 1 - 1/2 = 1/2$~~

~~$P(H) = 0.4 = p$~~

~~$P(T) = 1 - 0.4 = 0.6$~~

Expected value $\rightarrow E(x) = \sum_{i=1}^{\infty} x_i \cdot p(x_i)$

Mean, Variance and S.D

$x=0$	$/ x=1$
$P(x_0) = 0.4 = p = 1 - q$	$P(x=1) = 0.6$

$$= 0 \times 0.4 + 1 \times 0.6$$

$$= 0.6$$

formula for variance = $p(1-p) = pq$

" " STD & SD = \sqrt{pq}

Mean
P

Variance
 pq

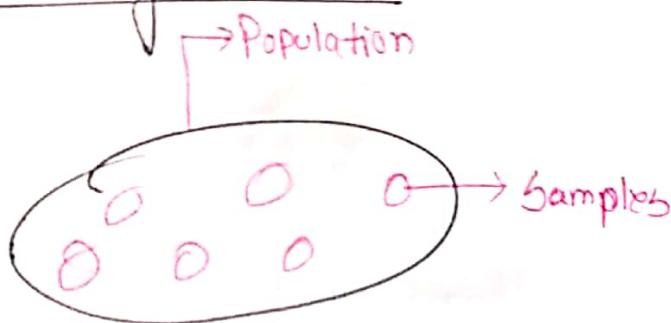
SD
 \sqrt{pq}

* If ratio is proportional, we can go toward random sampling better.

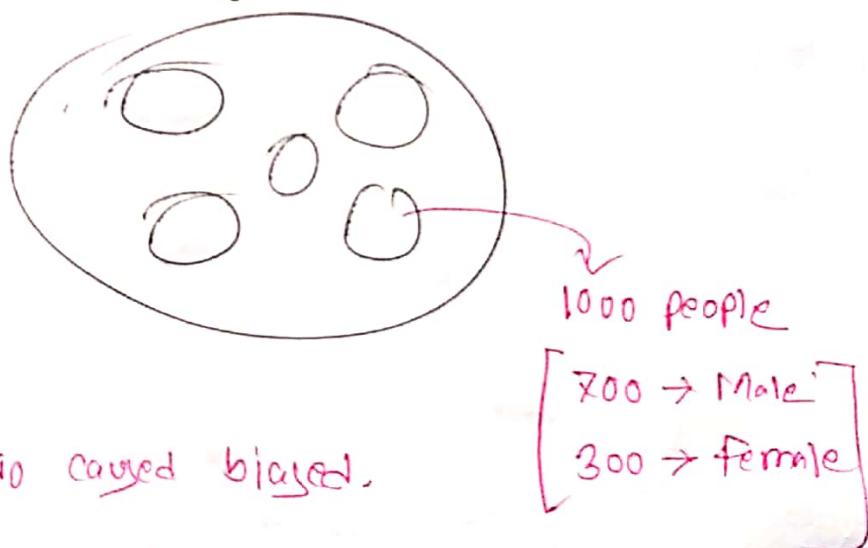
* For unbalanced class stratified Sampling technique works

Sampling Techniques

① Random Sampling Technique:



② Stratified Sampling technique:



7:3 ratio caused biased.

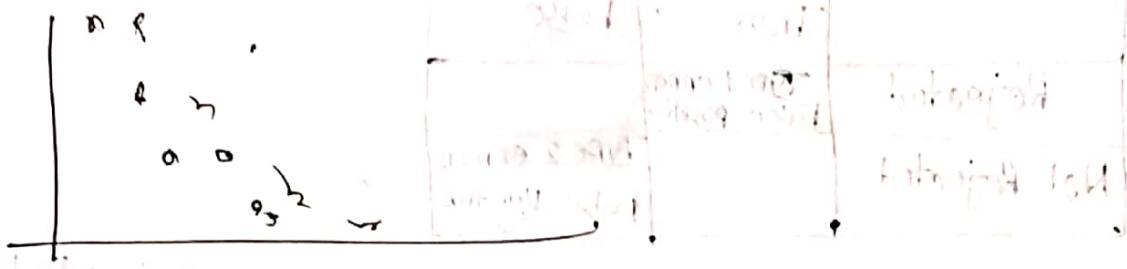
We've to provide 1:1 for good sampling.

③ Systematic Sampling:

Ex: If AT will take our job.
We will consider only domain expert



L1 and L2 regularization



math won = $\theta_0 + \theta_1 * \text{age}$ mathwon = $\theta_0 + \theta_1 * \text{age} + \theta_2 * \text{age}^2$
 $+ \theta_3 * \text{age}^3 + \theta_4 * \text{age}^4$

These regularization reduce overfitting by penalize higher degree of polynomial.

formula for L2 regularization $\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2 + \lambda \sum_{i=1}^n \theta_i^2$

Act like a knob. The bigger, theta value will be smaller.

$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - h_\theta(x_i))^2 + \lambda \sum_{i=1}^n |\theta_i|$

11

accuracy = $\frac{TP + TN}{P + N}$

Type 1 and Type 2 error

false positive / negative

Predicted

		P	N
Actual	P	TP	FN
	N	FP	TN

Type 1 error

Type 2 error

3 cases to consider - false positive / Negative:

- ① Whether a person has disease or not
- ② Market will crash or not
- ③ Vaccination side effect.

overdiagnosed patients (\rightarrow diagnosed not ①)

underdiagnosed patients (\rightarrow not diagnosed ②)

overdiagnosed patients (\rightarrow diagnosed not ③)

underdiagnosed patients (\rightarrow not diagnosed ④)

overdiagnosed patients (\rightarrow diagnosed not ⑤)

underdiagnosed patients (\rightarrow not diagnosed ⑥)

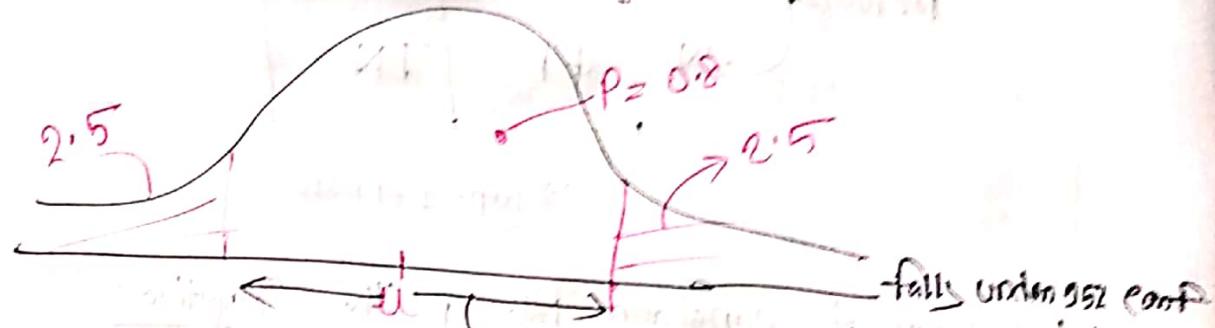
⑥ P-value also known as significance value, denoted by $\alpha = 0.05$

Hypothesis testing - Confidence, interval

Z-test / Statistics + p-value

Confidence interval: 68 - 95 - 99

Empirical formula: 95% confidence interval

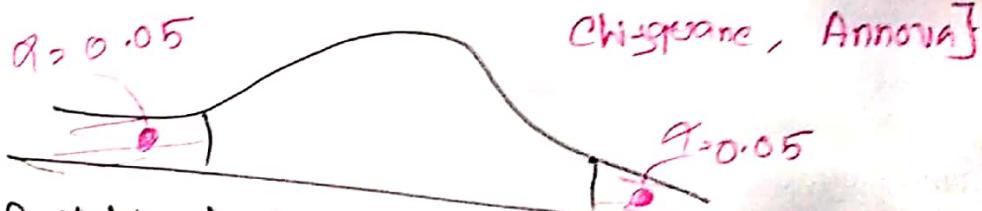


↳ $P = 0.95$ means for 100 experiments, that specific thing can happen 95 times.

Hypothesis testing - statistical

① Null hypothesis \rightarrow Alternative hypothesis

② Perform experiment $\rightarrow \alpha = 0.05$ { t-test, z-test, Chi-square, ANOVA}



If P value is less than α , then ~~there values~~ possibility to happen that will be rare.

So, we reject the null hypothesis.

③ Test value falls in extreme end \rightarrow Reject the null hypothesis.

Accept the alternative hypothesis.

Test value falls in confidence interval \rightarrow Accept null hypothesis.

Test value falls in confidence interval \rightarrow Accept null hypothesis.

Significance value \downarrow Confidence interval \uparrow Reject Alternative

If $\alpha = 0.05$ then CI is 95%.

For 2 tailed test

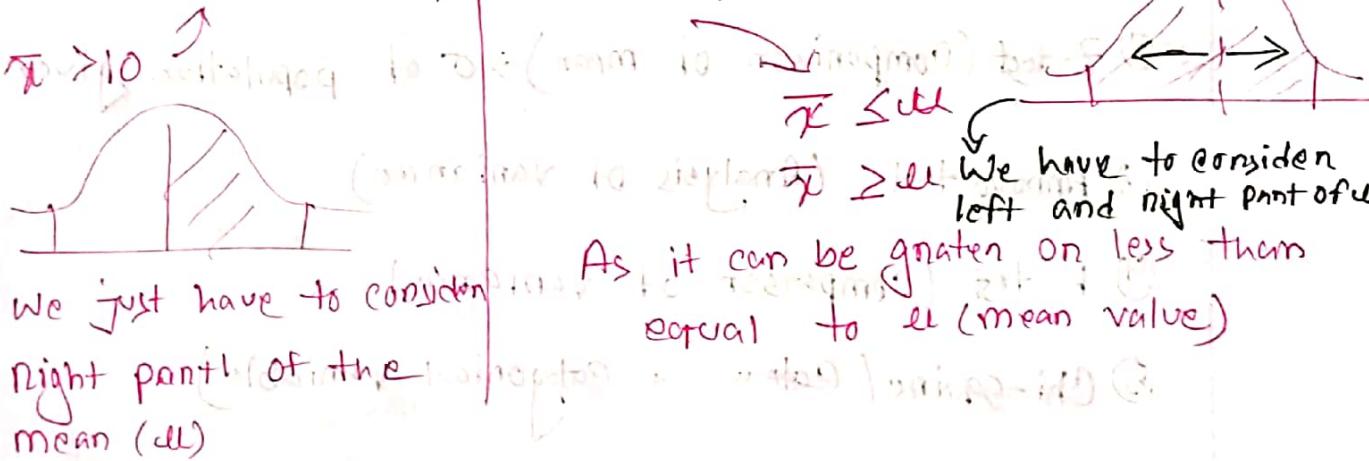
Experiment

2 tail test

$H_0 > 10$

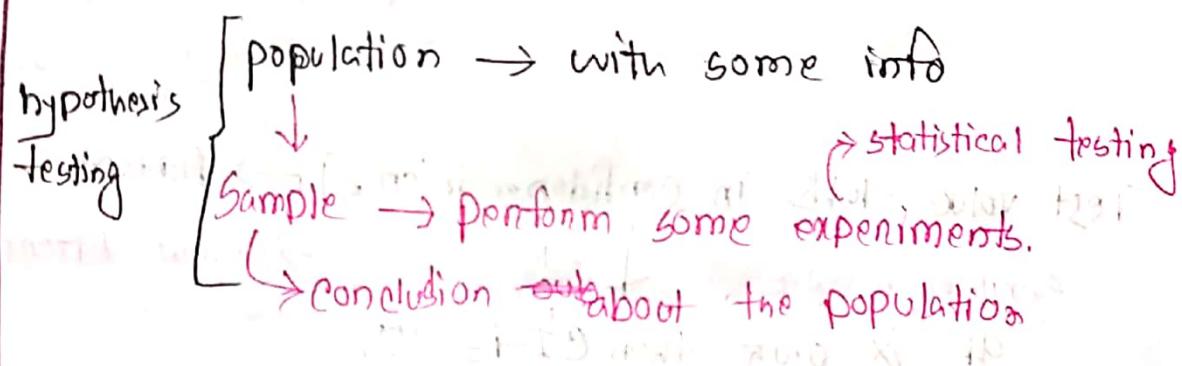
1 tail test

2 tail test



Case Study

Task: Avg height of people in India.



Statistical Test:

- ① T-test (Comparison of mean) → Sample (s) or given Standard deviation
- ② Z-test (Comparison of mean) → σ of population given
- ③ Anova-test (Analysis of variance)
- ④ F-test (Comparison of variance)
- ⑤ Chi-square (Categorical Variable)

z-test: What is the specific range in confidence interval.

(13)

not good

One Sample z-test



Next Page

IQ test for students:

Average, $\mu = 100$

Standard deviation, $\sigma = 15$

Doctor tested new medication to find out if it increases or decreases the IQ.

After 1 month ($n = 30$) were taken.

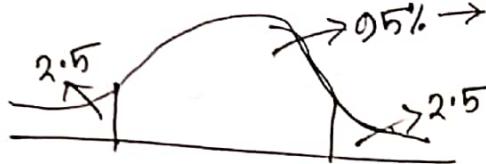
Where mean, $\bar{x} = 140$

Did medication affect intelligence. [Significance value; $\alpha = 0.05$]

①

⇒ Null hypothesis, $H_0 \Rightarrow \mu = 100$ } didn't affect.
so, this is a two-tail test

Alternative " , $H_1 \Rightarrow \mu \neq 100$



We need to find area under the curve of 0.9750
 $1 - \alpha = 0.9750$
 $1 - 0.025 = 0.9750$

② State the Alpha; Alpha is, $\alpha = 0.05$

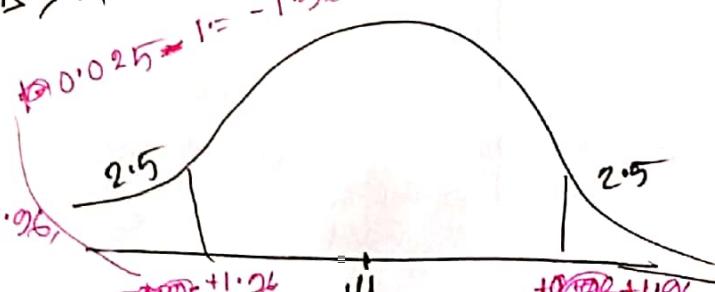
→ find specific region

③ State the decision rule

If they range between -1.96 to $+1.96$,

we accept the null hypothesis.

Otherwise ≤ -1.96 or $\geq +1.96$, we reject the null hypothesis and accept the alternative hypothesis.



\rightarrow different from
Z-score $\rightarrow \frac{1-5}{6}$

(4) Z-test:

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{140 - 100}{15 / \sqrt{30}}$$

$$= 14.60$$

(5) State Result:

Range to accept null hypothesis is -1.96 to $+1.96$

But Z_{score} is 14.60

As it doesn't fall into the range

\Rightarrow Reject the null hypothesis and accept the alternative

hypothesis.

AOC

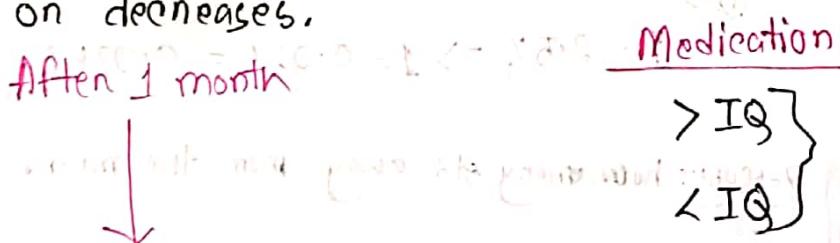
ROC

(M)

One sample Z test with application

In a population, the average IQ, $\mu = 100$ with $\sigma = 15$.

Doctor tested a new medication to find out whether it increases the iq or decreases.



Sample of 30 ($n=30$) participants were taken. Sample mean for 30 peoples IQ is 110 ($\bar{x} = 110$).

⇒ Did medication affect intelligence? Significance value is 0.05 ($\alpha = 0.05$)

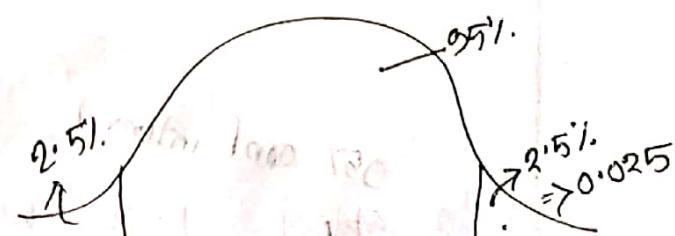
① ⇒ Hypothesis:

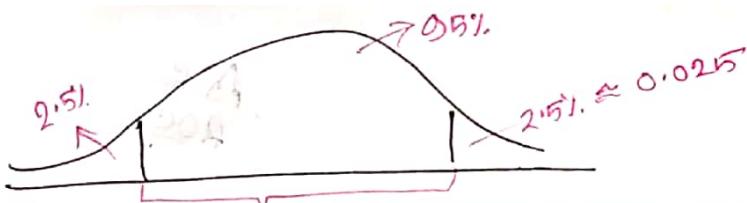
$H_0 \Rightarrow$ Did not affect the IQ ($\mu = 100$)

$H_1 \Rightarrow$ Affect the IQ ($\mu \neq 100$)

IQ can be increased or decreased. That's why it is a two tailed test.

② ⇒ State the alpha value: ($\alpha = 0.05$)





③ \Rightarrow State the decision rule: \leftarrow this range

Entire region for the above distribution is $= 1$

finding the range

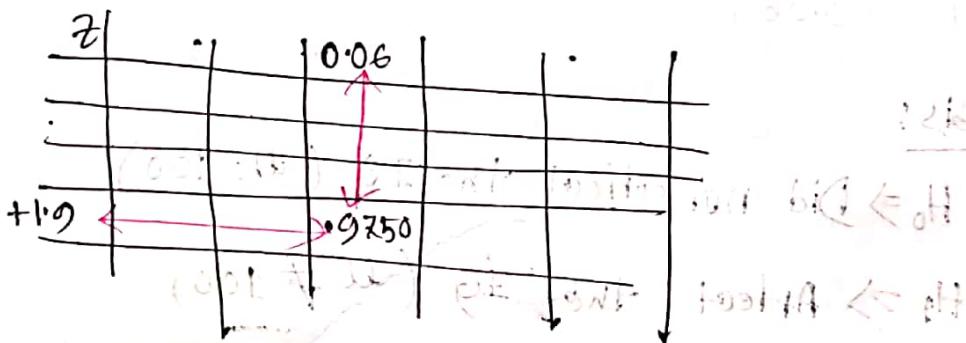
① \Rightarrow

$$1 - 2.5\% \Rightarrow 1 - 0.025 = 0.9750$$

Z-score: how many std away from the mean.

Z-test: find the range for 95% confidence interval.

② \Rightarrow Now we've to find out the area under the curve for 0.9750. We can find it on "Z-table".

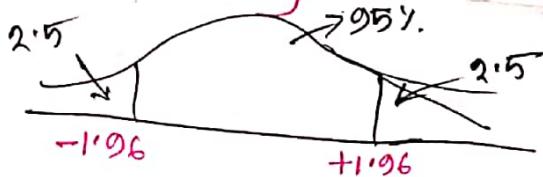


So, the area under the curve is: $1.9 + 0.06$

$$= 1.96$$

95% conf interval

\Rightarrow So, std is 1.96 away from the mean.



6

④ \Rightarrow Z-test: different from Z-score

$$\text{formula: } z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{140 - 100}{\frac{15}{\sqrt{30}}} = 14.60$$

⑤ State the negt: We'll accept the null hypothesis if z-test

Value fall into the area under the curve range. (-1σ to +1σ)

14.60 doesn't fall into -1.96 to +1.96

\Rightarrow Reject the null hypothesis.]

\Rightarrow Accept \neg Alternative

Changes in IG

(*) One sample z-test with proportion:

A survey claim that 9 out of 10 doctors recommended aspirin for their patient with headache. To test this claim, a random sample of 100 doctors is taken. Out of this 100 doctors, 82 indicate that they recommend Aspirin. $\alpha = 0.05$

① \Rightarrow Hypothesis:

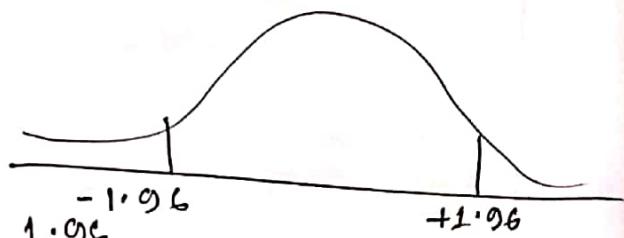
Null Hypothesis, $H_0 = P = 0.90$

Alternative " " " $H_1 = P \neq 0.90$

② \Rightarrow Significance: $\alpha = 0.05$

Range: -1.96 to +1.96

As $1 - 0.05 = 0.95 \rightarrow$ z-table $\rightarrow 1.96$



③ \Rightarrow Decision Rule:

④ \Rightarrow Z-statistics:

$$\text{formula, } Z_0 = \frac{\hat{P} - P_0}{\sqrt{\frac{P_0(1-P_0)}{n}}}$$
$$= \frac{0.82 - 0.90}{\sqrt{\frac{0.90(0.10)}{100}}} = -2.66$$

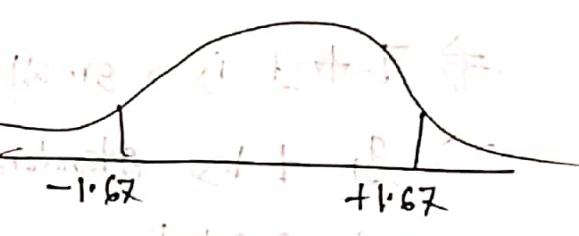
Probability respect to sample
 $H_0, P = 0.82 \Rightarrow \frac{82}{100}$
 $P_0 = 0.90 \Rightarrow$ Probability respect to population
 $n = 100 \Rightarrow \frac{9}{10}$
num of samples

(b) State the result:

Z value is $= -2.667$

Our range is $\downarrow -1.67$ to $+1.67$

from Z table



As Z value is out of the range, we reject the null hypothesis.

∴

Test H₀ for no significant difference is rejected.

∴ H_1 is accepted.

∴ H_0 is rejected.

mid sample

H_1 is accepted.

∴ H_0 is rejected.

∴ H_1 is accepted.

T-test

- ⇒ T-test is a small sample test.
- ⇒ In t-test standard deviation (σ) is given for sample, not population.

⇒ Formula:

$$t = \frac{\bar{x} - \mu_0}{\frac{\sigma \rightarrow \text{sample}}{\sqrt{n}}}$$

⇒ App

Application of t-test

- ① Size of sample is small ($n < 30$)
- ② Degree of freedom is ($v = n - 1$) \rightarrow Sample size
- ⇒ T-test is used for test of significance of regression coefficient in regression model.

Statistical testing

Parametric test

Non-parametric test

① Parametric test: Suitable for normally distributed data.

② Non-Parametric test: Suitable for any continuous data, based on ranks of the data values.

⇒ Three types of t-test:

① One sample t-test

② Unpaired, two sample t-test

③ Paired sample t-test.

⇒ Conditions for t-test:

① The data has to follow a continuous or ordinal scale.

② The data has to be randomly selected.

③ The data should be normally distributed.

⇒ t-score:

The t-score formula enable us to transform a distribution into a standardize form, which we use to compare the score.

$$\text{formula, } t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{N_1} + \frac{s_2^2}{N_2}}}$$

s_1^2, s_2^2 : Varience

N_1, N_2 : Num of samples

\bar{x}_1, \bar{x}_2 : mean of two groups

Group	Method	n	\bar{x}	s
1	Intensive	12	46.31	6.4
2	Paced	10	42.79	7.52

$$t = \frac{46.31 - 42.79}{\sqrt{\frac{6.4^2}{12} + \frac{7.52^2}{10}}} = \frac{3.52}{\sqrt{3.46 + 5.66}} = 1.166$$

t value is $= 1.166$. Now to find valid range, we need to look into t-table with significance value and degree of freedom.

degree of freedom is smallest among $12-1, 10-1$

∴ Degree of freedom = 9

$$\alpha = 0.05$$

With df = 9 and $\alpha = 0.05$, from t-table, we

find that critical value is $= 1.833$

As t-value 1.66 doesn't exceed 1.833 , we accept the null hypothesis.

No std for population
std for sample.

One Sample t-test

Z-tests are commonly done when population std is known

t-tests " " " " " " is unknown

* We use t-distribution as population std is unknown.
that we have to estimate.

→ Example Avg IQ 100. Scientist experimenting that medication has positive/negative effect.

To test affect of the medication, 30 participants were taken with mean of 140 and std of 20. Alpha = 0.05

① Define Hypothesis:

$$H_0 \Rightarrow \mu = 100$$

$$H_1 \Rightarrow \mu \neq 100$$

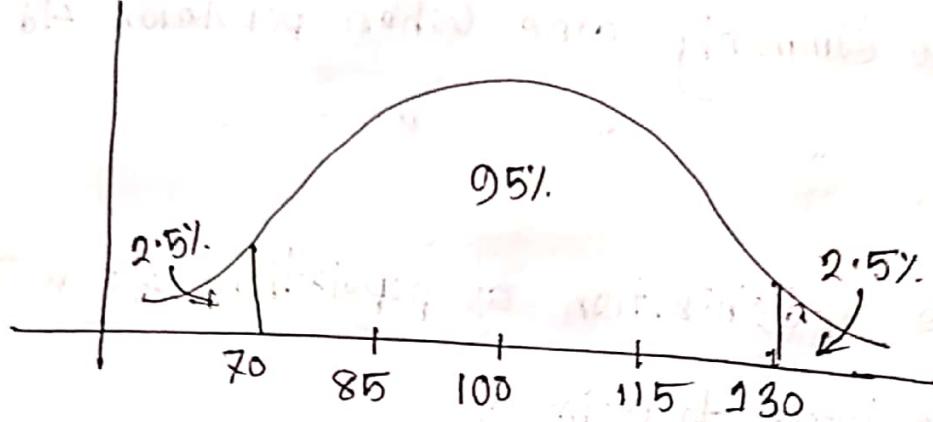
② State alpha:

$$\alpha = 0.05$$

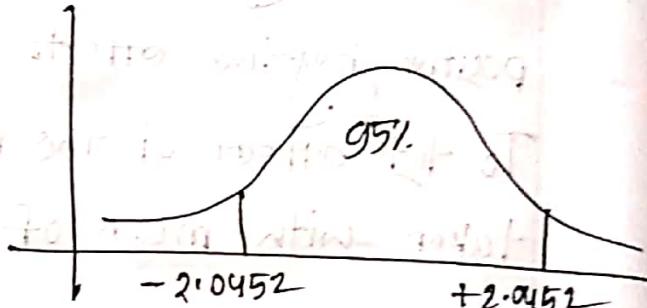
③ Calculate Degree of freedom:

$$n - 1 = 30 - 1 = 29; n = \text{num of sample.}$$

④ \Rightarrow State Decision Rule!



df		Critical Value
2tailed test	0.1	0.5
29		2.0952



⇒ Critical value is = -2.0452 to +2.0452

if t is greater/less than critical value, reject the null hypothesis.

⑤ \Rightarrow Calculate test statistics:

$$t = \frac{\bar{x} - \mu}{\sigma}$$

$$= \frac{140 - 100}{20}$$

$$= 10.96$$

$$\bar{x} = 140$$

III = 190

$$S = 20$$

$$n = 30$$

(P)

6) State result:

Critical value is: $-2.0452 - 2.0452$

$$t = 10.96$$

As t value is far away from critical value, reject the null hypothesis.

if changed the TG .

Independent Sample t-test

⇒ Compute two independent sample

⇒ Comp two class. Class A: 25 student, $\bar{u}_1 = 70, s_1 = 15$
 Class B: 20 student, $\bar{u}_2 = 79, s_2 = 25$. $\alpha = 0.05$. Are they different

$$H_0 \Rightarrow \bar{u}_1 = \bar{u}_2$$

$$H_1 \Rightarrow \bar{u}_1 \neq \bar{u}_2$$

$$\Rightarrow \alpha = 0.05$$

$$\Rightarrow \text{Degree of freedom} \\ df = (n_1 - 1) + (n_2 - 1) = (25 - 1) + (20 - 1) = 43$$

2) Decision rule: Critical range from t-table

$$\rightarrow -2.0167 \longleftrightarrow 2.0167$$

⑦ \Rightarrow test stat:

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}}$$

$$s_p^2 = \frac{ss_1 + ss_2}{df_1 + df_2}$$

$$\begin{aligned} &= \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_p^2}{n_1} + \frac{s_p^2}{n_2}}} \\ &= \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{401.74}{25} + \frac{401.74}{20}}} \\ &= \frac{-4}{\sqrt{36.16}} \end{aligned}$$

$$\begin{aligned} &= \frac{5400 + 11875}{29 + 19} \\ &= 401.74 \end{aligned}$$

$$\begin{aligned} df_1 &= n-1 = 25-1 = 24 \\ df_2 &= n-1 = 20-1 = 19 \\ ss_1 &= s_1^2(df_1) = 15^2(24) = 5400 \\ ss_2 &= (25^2)(10) = 11875 \\ s_2^2 &= 11875 \end{aligned}$$

$$t = -0.67$$

⑧ \Rightarrow state resol:

Decision Rule: Critical range $-2.0157 \leftarrow 2.0157$

$$t = -0.67$$

Accept the null hypothesis

$t < -2.0157$

$0 < t < 2.0157$

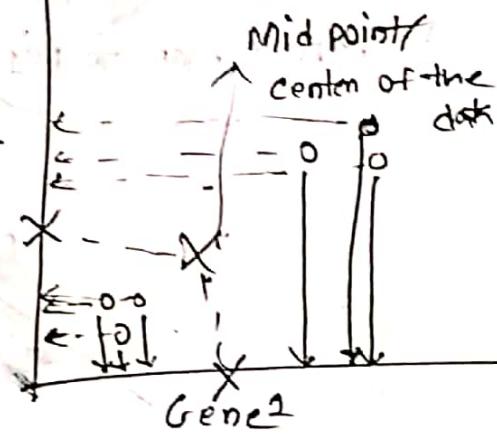
Conclusion: H_0 is not rejected (H_0 is true) \rightarrow no difference between sample

Tested null: H_0 is true (H_0 is not rejected)

$> 20.5 \rightarrow$ false $\rightarrow H_0$ is not rejected

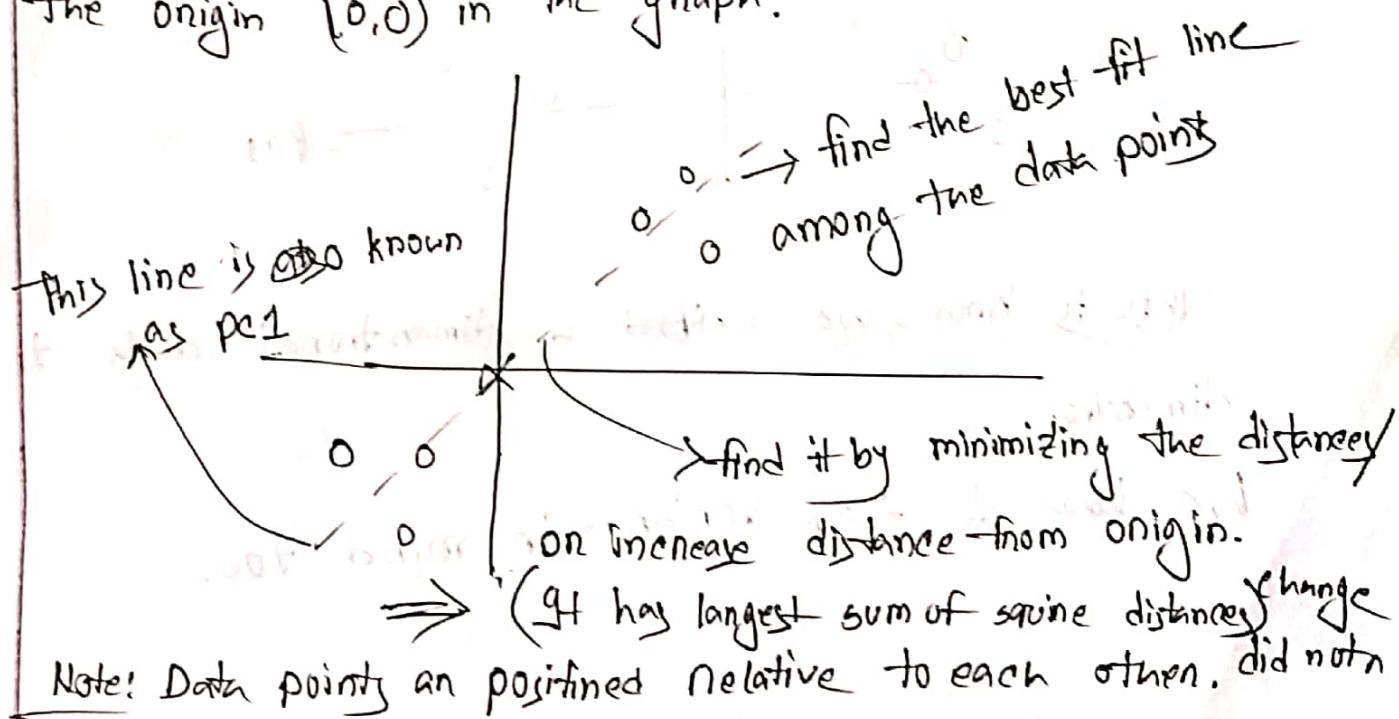
PCA

	Mouse 1	M2	M3	M4	M5	M6
Gene 1	10	11	8	3	2	1
Gene 2	6	4	5	3	2.8	1

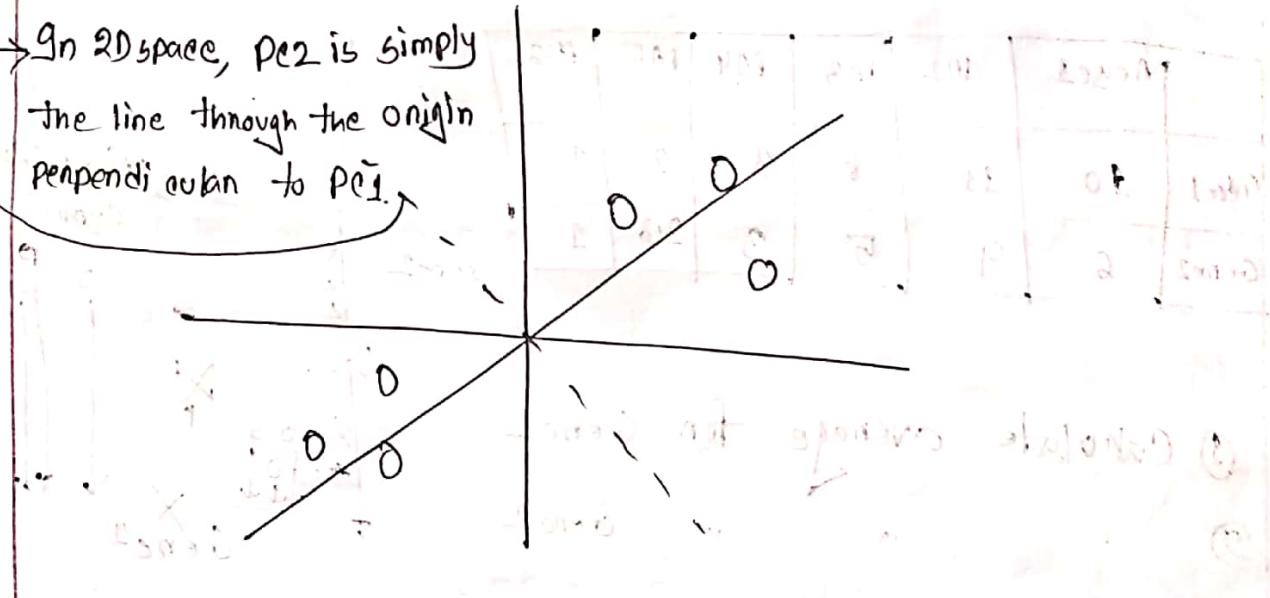


- ① Calculate average for Gene 1
- ② " " " " Gene 2

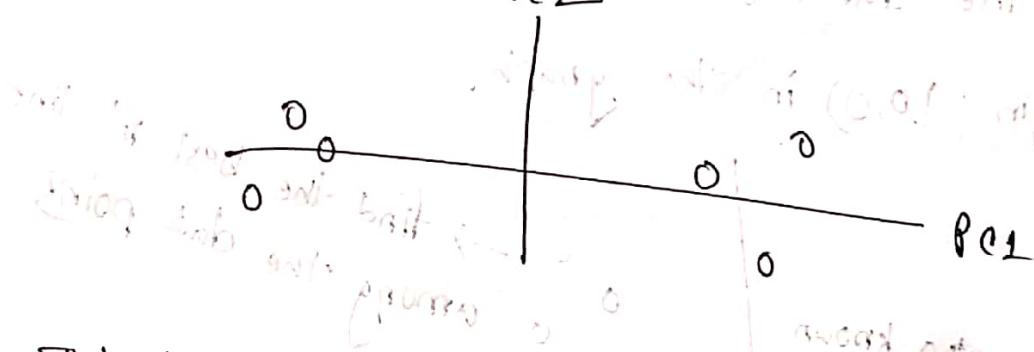
* Shift the data, so that the center is on top of the origin (0,0) in the graph.



In 2D space, PC_2 is simply the line through the origin perpendicular to PC_1 .



finally shift them to horizontal.



This is how we shifted 6 dimensional data to two dimension.

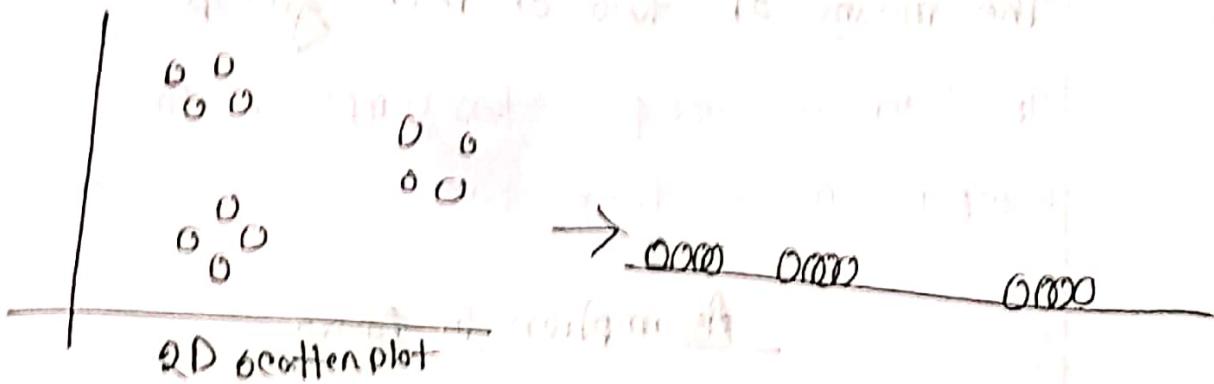
We have to keep variance in mind too.

most axes are aligned (12) ←
most axes are aligned (12) ←

t-5NE

93

⇒ It takes high dimensional data. Reduce it to low dimensional graph that retain a lot of original information.



Let's see how 1-SHR transforms this graph into a flat, 1-D plot on a number line.

t-SNE finds a way to project data into low dimensional space, so that the clustering in the high dimensional space is preserved.

Step 1: put the points on the number plot in a random order.

0 0 0 0 0 0 0 0 0 180916 10:16:23

- SNE move them a bit by bit until it cluster them

→ q.t. originated by simikh points and pushed by not simikh once

gl moves according to the orthoaction / ditraction

Anova test

Anova \Rightarrow Analysis of variance

Anova is a statistical method used to compare the means of two or more groups.

In t-test, we compare two sample, in anova, we can compare more than two.

Assumption in Anova

1. Normality of sampling distribution of mean.

② Independence of error

③ Absence of outliers.

④ Homogeneity of variance

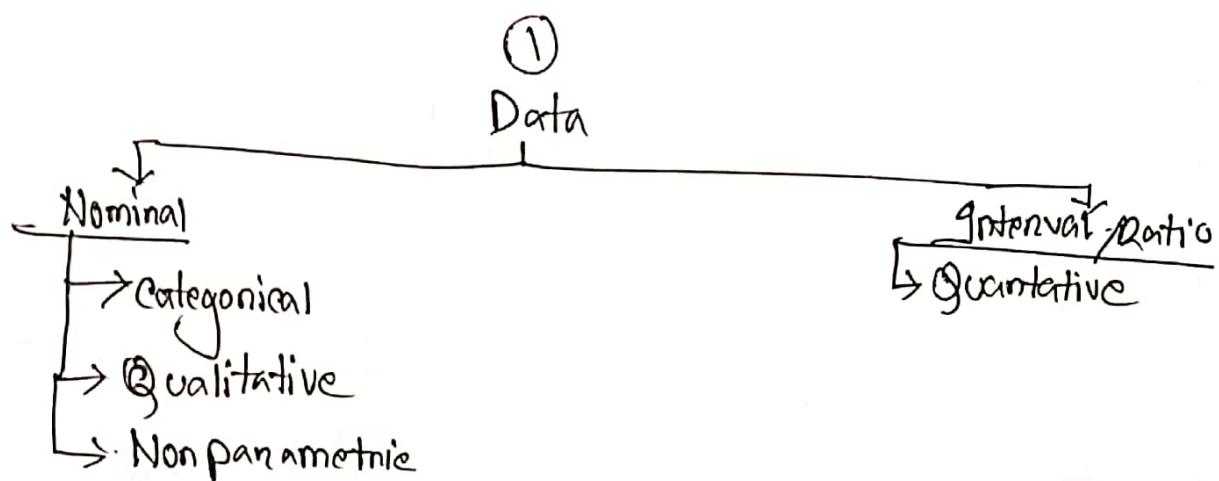
Hypothesis

$$H_0 \Rightarrow \mu_{A_1} = \mu_{A_2} = \mu_{A_3}$$

$H_1 \Rightarrow$ Not all means are equal.

Choosing the test

- ① Data
- ② Sample
- ③ Purpose.



Test statistics for nominal data:

- ① Test for proportion
- ② Difference of two proportion
- ③ Chi-square test for independence

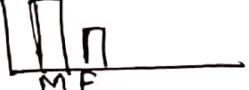
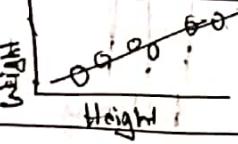
Test statistics for interval data:

- ① Test for a mean
- ② Difference of two means (Paired)
- ③ Difference of two means (Independent Samples)
- ④ Regression Analysis.

What to use, where

t-test, Chi-square test, p value

Gender	Age	Height	Weight
Female	Adult	1.9	60
Male	Child	1.2	15
Male	Adult	1.5	85
Female	"	1.3	74
Male	"	1.6	72
Female	Elderly	1.5	65

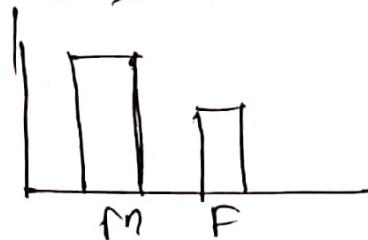
	What we observe in our sample data	As it real
One categorical		1 sample proportion test
Two "		Chi squared
One numerical		t-test
One numerical & One categorical		t-test on Anova
two numerical		Correlation test

① \Rightarrow One category: (Gender)

$H_0 \Rightarrow$ There is no difference in gender proportion



$H_1 \Rightarrow$ There is a difference:



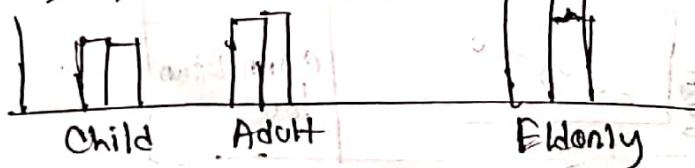
Gender	Males	Females	Total
Male	16,67	13,33	30,00
Female	8,33	11,67	20,00
Total	25	25	50

Test! We need to know how significant this difference is by One sample proportion test.

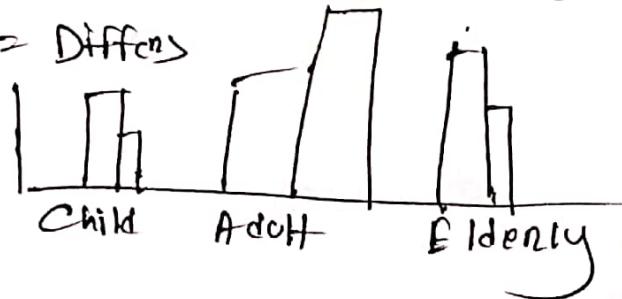
② \Rightarrow Two Categories: (Gender, Age)

Does the proportion of males and females differ across age groups?

$H_0 =$ Does not diffens.



$H_1 =$ Diffens



⇒ test: chi-square → calculate p-value

if p is less than α , reject null hypothesis

② One numeric! (height)

Is height different from previous year?

$$\frac{8}{17} \cdot (307) = 17.6 \quad \frac{9}{17} \cdot (307) = 17.1$$

total food and meat portion of 1st class passengers are 12.8 kg
with first class more likely to contribute 30% to
first class food and meat portion of 1st class passengers are 12.8 kg

first class lower limit

$$17.6 - 1.7 = 15.9 \quad \frac{8}{17} \cdot (307 + 1.7) = 16.9$$

$$17.1 - 1.7 = 15.4 \quad 17.1 + 1.7 = 18.8 \quad 17.1 \times (307 + 1.7) = 18.8$$

$$\frac{8}{17} \times \frac{5}{17} = \frac{40}{289}$$

$$\frac{9}{17} \times \frac{5}{17} = \frac{45}{289}$$

$$\frac{8}{17} \times \frac{9}{17} = \frac{72}{289}$$

$$\frac{9}{17} \times \frac{9}{17} = \frac{81}{289}$$

$$\frac{8}{17} \times \frac{8}{17} = \frac{64}{289}$$

$$\frac{9}{17} \times \frac{8}{17} = \frac{72}{289}$$

$$\frac{8}{17} \times \frac{7}{17} = \frac{56}{289}$$

$$\frac{9}{17} \times \frac{7}{17} = \frac{63}{289}$$

$$\frac{8}{17} \times \frac{6}{17} = \frac{48}{289}$$

$$\frac{9}{17} \times \frac{6}{17} = \frac{54}{289}$$

$$\frac{8}{17} \times \frac{5}{17} = \frac{40}{289}$$

$$\frac{9}{17} \times \frac{5}{17} = \frac{45}{289}$$

$$\frac{8}{17} \times \frac{4}{17} = \frac{32}{289}$$

$$\frac{9}{17} \times \frac{4}{17} = \frac{36}{289}$$

$$\frac{8}{17} \times \frac{3}{17} = \frac{24}{289}$$

$$\frac{9}{17} \times \frac{3}{17} = \frac{27}{289}$$

$$\frac{8}{17} \times \frac{2}{17} = \frac{16}{289}$$

$$\frac{9}{17} \times \frac{2}{17} = \frac{18}{289}$$

$$\frac{8}{17} \times \frac{1}{17} = \frac{8}{289}$$

$$\frac{9}{17} \times \frac{1}{17} = \frac{9}{289}$$

$$\frac{8}{17} \times \frac{0}{17} = \frac{0}{289}$$

$$\frac{9}{17} \times \frac{0}{17} = \frac{0}{289}$$

$$\frac{8}{17} \times \frac{1}{17} = \frac{8}{289}$$

$$\frac{9}{17} \times \frac{1}{17} = \frac{9}{289}$$

$$\frac{8}{17} \times \frac{2}{17} = \frac{16}{289}$$

$$\frac{9}{17} \times \frac{2}{17} = \frac{18}{289}$$

$$\frac{8}{17} \times \frac{3}{17} = \frac{24}{289}$$

$$\frac{9}{17} \times \frac{3}{17} = \frac{27}{289}$$

$$\frac{8}{17} \times \frac{4}{17} = \frac{32}{289}$$

$$\frac{9}{17} \times \frac{4}{17} = \frac{36}{289}$$

$$\frac{8}{17} \times \frac{5}{17} = \frac{40}{289}$$

$$\frac{9}{17} \times \frac{5}{17} = \frac{45}{289}$$

$$\frac{8}{17} \times \frac{6}{17} = \frac{48}{289}$$

$$\frac{9}{17} \times \frac{6}{17} = \frac{54}{289}$$

$$\frac{8}{17} \times \frac{7}{17} = \frac{56}{289}$$

$$\frac{9}{17} \times \frac{7}{17} = \frac{63}{289}$$

$$\frac{8}{17} \times \frac{8}{17} = \frac{64}{289}$$

$$\frac{9}{17} \times \frac{8}{17} = \frac{72}{289}$$

$$\frac{8}{17} \times \frac{9}{17} = \frac{72}{289}$$

$$\frac{9}{17} \times \frac{9}{17} = \frac{81}{289}$$

$$\frac{8}{17} \times \frac{10}{17} = \frac{80}{289}$$

$$\frac{9}{17} \times \frac{10}{17} = \frac{90}{289}$$

$$\frac{8}{17} \times \frac{11}{17} = \frac{88}{289}$$

$$\frac{9}{17} \times \frac{11}{17} = \frac{99}{289}$$

$$\frac{8}{17} \times \frac{12}{17} = \frac{96}{289}$$

$$\frac{9}{17} \times \frac{12}{17} = \frac{108}{289}$$

$$\frac{8}{17} \times \frac{13}{17} = \frac{104}{289}$$

$$\frac{9}{17} \times \frac{13}{17} = \frac{117}{289}$$

$$\frac{8}{17} \times \frac{14}{17} = \frac{112}{289}$$

$$\frac{9}{17} \times \frac{14}{17} = \frac{126}{289}$$

$$\frac{8}{17} \times \frac{15}{17} = \frac{120}{289}$$

$$\frac{9}{17} \times \frac{15}{17} = \frac{135}{289}$$

$$\frac{8}{17} \times \frac{16}{17} = \frac{128}{289}$$

$$\frac{9}{17} \times \frac{16}{17} = \frac{144}{289}$$

$$\frac{8}{17} \times \frac{17}{17} = \frac{136}{289}$$

$$\frac{9}{17} \times \frac{17}{17} = \frac{153}{289}$$

$$\frac{8}{17} \times \frac{18}{17} = \frac{144}{289}$$

$$\frac{9}{17} \times \frac{18}{17} = \frac{162}{289}$$

$$\frac{8}{17} \times \frac{19}{17} = \frac{152}{289}$$

$$\frac{9}{17} \times \frac{19}{17} = \frac{171}{289}$$

$$\frac{8}{17} \times \frac{20}{17} = \frac{160}{289}$$

$$\frac{9}{17} \times \frac{20}{17} = \frac{180}{289}$$

$$\frac{8}{17} \times \frac{21}{17} = \frac{168}{289}$$

$$\frac{9}{17} \times \frac{21}{17} = \frac{189}{289}$$

$$\frac{8}{17} \times \frac{22}{17} = \frac{176}{289}$$

$$\frac{9}{17} \times \frac{22}{17} = \frac{198}{289}$$

$$\frac{8}{17} \times \frac{23}{17} = \frac{184}{289}$$

$$\frac{9}{17} \times \frac{23}{17} = \frac{207}{289}$$

$$\frac{8}{17} \times \frac{24}{17} = \frac{192}{289}$$

$$\frac{9}{17} \times \frac{24}{17} = \frac{216}{289}$$

$$\frac{8}{17} \times \frac{25}{17} = \frac{200}{289}$$

$$\frac{9}{17} \times \frac{25}{17} = \frac{225}{289}$$

$$\frac{8}{17} \times \frac{26}{17} = \frac{208}{289}$$

$$\frac{9}{17} \times \frac{26}{17} = \frac{234}{289}$$

$$\frac{8}{17} \times \frac{27}{17} = \frac{216}{289}$$

$$\frac{9}{17} \times \frac{27}{17} = \frac{243}{289}$$

$$\frac{8}{17} \times \frac{28}{17} = \frac{224}{289}$$

$$\frac{9}{17} \times \frac{28}{17} = \frac{252}{289}$$

$$\frac{8}{17} \times \frac{29}{17} = \frac{232}{289}$$

$$\frac{9}{17} \times \frac{29}{17} = \frac{261}{289}$$

$$\frac{8}{17} \times \frac{30}{17} = \frac{240}{289}$$

$$\frac{9}{17} \times \frac{30}{17} = \frac{270}{289}$$

$$\frac{8}{17} \times \frac{31}{17} = \frac{248}{289}$$

$$\frac{9}{17} \times \frac{31}{17} = \frac{279}{289}$$

$$\frac{8}{17} \times \frac{32}{17} = \frac{256}{289}$$

$$\frac{9}{17} \times \frac{32}{17} = \frac{288}{289}$$

$$\frac{8}{17} \times \frac{33}{17} = \frac{264}{289}$$

$$\frac{9}{17} \times \frac{33}{17} = \frac{297}{289}$$

$$\frac{8}{17} \times \frac{34}{17} = \frac{272}{289}$$

$$\frac{9}{17} \times \frac{34}{17} = \frac{306}{289}$$

$$\frac{8}{17} \times \frac{35}{17} = \frac{280}{289}$$

$$\frac{9}{17} \times \frac{35}{17} = \frac{315}{289}$$

$$\frac{8}{17} \times \frac{36}{17} = \frac{288}{289}$$

$$\frac{9}{17} \times \frac{36}{17} = \frac{324}{289}$$

$$\frac{8}{17} \times \frac{37}{17} = \frac{296}{289}$$

$$\frac{9}{17} \times \frac{37}{17} = \frac{333}{289}$$

$$\frac{8}{17} \times \frac{38}{17} = \frac{304}{289}$$

$$\frac{9}{17} \times \frac{38}{17} = \frac{342}{289}$$

$$\frac{8}{17} \times \frac{39}{17} = \frac{312}{289}$$

$$\frac{9}{17} \times \frac{39}{17} = \frac{351}{289}$$

$$\frac{8}{17} \times \frac{40}{17} = \frac{320}{289}$$

$$\frac{9}{17} \times \frac{40}{17} = \frac{360}{289}$$

$$\frac{8}{17} \times \frac{41}{17} = \frac{328}{289}$$

$$\frac{9}{17} \times \frac{41}{17} = \frac{369}{289}$$

$$\frac{8}{17} \times \frac{42}{17} = \frac{336}{289}$$

$$\frac{9}{17} \times \frac{42}{17} = \frac{378}{289}$$

$$\frac{8}{17} \times \frac{43}{17} = \frac{344}{289}$$

$$\frac{9}{17} \times \frac{43}{17} = \frac{387}{289}$$

$$\frac{8}{17} \times \frac{44}{17} = \frac{352}{289}$$

$$\frac{9}{17} \times \frac{44}{17} = \frac{396}{289}$$

$$\frac{8}{17} \times \frac{45}{17} = \frac{360}{289}$$

$$\frac{9}{17} \times \frac{45}{17} = \frac{405}{289}$$

$$\frac{8}{17} \times \frac{46}{17} = \frac{368}{289}$$

$$\frac{9}{17} \times \frac{46}{17} = \frac{414}{289}$$

$$\frac{8}{17} \times \frac{47}{17} = \frac{376}{289}$$

$$\frac{9}{17} \times \frac{47}{17} = \frac{423}{289}$$

$$\frac{8}{17} \times \frac{48}{17} = \frac{384}{289}$$

$$\frac{9}{17} \times \frac{48}{17} = \frac{432}{289}$$

$$\frac{8}{17} \times \frac{49}{17} = \frac{392}{289}$$

$$\frac{9}{17} \times \frac{49}{17} = \frac{441}{289}$$

$$\frac{8}{17} \times \frac{50}{17} = \frac{400}{289}$$

$$\frac{9}{17} \times \frac{50}{17} = \frac{450}{289}$$

$$\frac{8}{17} \times \frac{51}{17} = \frac{408}{289}$$

$$\frac{9}{17} \times \frac{51}{17} = \frac{459}{289}$$

$$\frac{8}{17} \times \frac{52}{17} = \frac{416}{289}$$

$$\frac{9}{17} \times \frac{52}{17} = \frac{468}{289}$$

$$\frac{8}{17} \times \frac{53}{17} = \frac{424}{289}$$

$$\frac{9}{17} \times \frac{53}{17} = \frac{477}{289}$$

$$\frac{8}{17} \times \frac{54}{17} = \frac{432}{289}$$

$$\frac{9}{17} \times \frac{54}{17} = \frac{486}{289}$$

$$\frac{8}{17} \times \frac{55}{17} = \frac{440}{289}$$

$$\frac{9}{17} \times \frac{55}{17} = \frac{495}{289}$$

$$\frac{8}{17} \times \frac{56}{17} = \frac{448}{289}$$

$$\frac{9}{17} \times \frac{56}{17} = \frac{504}{289}$$

$$\frac{8}{17} \times \frac{57}{17} = \frac{456}{289}$$

$$\frac{9}{17} \times \frac{57}{17} = \frac{513}{289}$$

$$\frac{8}{17} \times \frac{58}{17} = \frac{464}{289}$$

$$\frac{9}{17} \times \frac{58}{17} = \frac{522}{289}$$

$$\frac{8}{17} \times \frac{59}{17} = \frac{472}{289}$$

$$\frac{9}{17} \times \frac{59}{17} = \frac{531}{289}$$

$$\frac{8}{17} \times \frac{60}{17} = \frac{480}{289}$$

$$\frac{9}{17} \times \frac{60}{17} = \frac{540}{289}$$

$$\frac{8}{17} \times \frac{61}{17} = \frac{488}{289}$$

$$\frac{9}{17} \times \frac{61}{17} = \frac{549}{289}$$

$$\frac{8}{17} \times \frac{62}{17} = \frac{496}{289}$$

$$\frac{9}{17} \times \frac{62}{17} = \frac{558}{289}$$

$$\frac{8}{17} \times \frac{63}{17} = \frac{504}{289}$$

$$\frac{9}{17} \times \frac{63}{17} = \frac{567}{289}$$

$$\frac{8}{17} \times \frac{64}{17} = \frac{512}{289}$$

$$\frac{9}{17} \times \frac{64}{17} = \frac{576}{289}$$

$$\frac{8}{17} \times \frac{65}{17} = \frac{520}{289}$$

$$\frac{9}{17} \times \frac{65}{17} = \frac{585}{289}$$

$$\frac{8}{17} \times \frac{66}{17} = \frac{528}{289}$$

Bayes theorem

① Conditional probability

② Independent event → coin toss → doesn't depend on each other

③ Dependent event → one event affecting the probability of another event

$$P(\text{Black}) = \frac{3}{10}$$

0	0	0	0	0	0	0
0	0	0				

$$P(\text{Pink}) = \frac{7}{10}$$

dependent event

If we randomly select two marbel from this box, what is the probability of drawing a ~~green~~^{pink} marbel and then a ~~blue~~^{black} marbel, without replacement.

first event

$$\begin{aligned} P(\text{pink marbel}) &= \frac{7}{10} \\ &= 0.7 \end{aligned}$$

Second event:

$$\begin{aligned} P(\text{black marbel}) &= \frac{3}{7} \rightarrow \text{as we drawn} \\ &= 0.33 \quad \text{one marbel} \end{aligned}$$

$$\begin{aligned} P(A \cap B) &= P(A) \times P(B) \\ &= \frac{7}{10} \times \frac{3}{7} \\ &= \frac{7}{30} \\ &= 0.233 \end{aligned}$$

Conditional probability,

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{n(A \cap B)}{n(B)}$$

* Probability of Given that / probability of picking a pink manbel is $\frac{7}{10}$. What is the probability of picking black manbel given for given probability of pink manbel.

	पुरुष	महिला	
फ्रांसीसी	50	90	140
अंग्रेजी	90	60	150
सामग्री	140	150	290

① पुरुष युवा फ्रांसीसी का संख्यावाला मौका ?

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{50}{140} = \frac{5}{14}$$

Bayes theorem

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$\Rightarrow P(A \cap B) = P(B \cap A)$$

$$\rightarrow P(A \cap B) = P(A|B) * P(B)$$

$$\rightarrow P(B \cap A) = P(B|A) * P(A)$$

$$\text{As } P(A \cap B) = P(B \cap A)$$

$$= P(A|B) * P(B) = P(B|A) * P(A)$$

$$\therefore P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Posterior
Probability

Marginal Probability

Naive Bayes classification

$$\text{Bayes theorem, } P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

Dataset

$$\{x_1, x_2, x_3, x_4, \dots, x_n\} \quad \{y\}$$

$$P(y|x_1, x_2, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)P(x_3|y)\dots P(x_n|y) * P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$= \frac{P_y * \prod_{i=1}^n P(x_i|y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$\therefore P(y|x_1, x_2, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

$$y = \arg \max_y P(y) \prod_{i=1}^n P(x_i|y)$$

→ We'll accept the highest value for that classifier

Suppose for y happy

$$Yes = 0.7$$

$$No = 0.3$$

⇒ We'll accept 0.7 → Yes

Outlook

	Yes	No	P(Y)	P(No)
Sunny	2	3	2/5	3/5
Overcast	9	0	9/9	0/9
Rainy	3	2	3/5	2/5
Total	9	5	100%	100%

Temperature

	Yes	No	P(Y)	P(No)
HOT	2	2	2/5	2/5
MILD	4	2	4/9	2/9
COLD	3	1	3/5	2/5
Total	9	5	100%	100%

→ Suppose today is sunny and Hot. What is the probability of ~~training~~, playing?

Today (sunny, hot)

	P(yes)	P(yes)
Yes	0	9/14
No	5	5/14
Total	14	

$$P(Y \mid \text{sunny, hot}) = \frac{P(\text{sunny} \mid \text{yes}) * P(\text{hot} \mid \text{yes}) * P(\text{yes})}{P(\text{hot}) * P(\text{sunny})} \rightarrow \text{ignore}$$

$$\begin{aligned} P(Y \mid \text{sunny, hot}) &= \frac{2}{9} * \frac{2}{5} * \frac{9}{14} \\ &= 0.031 \end{aligned}$$

$$\begin{aligned} P(N \mid \text{sunny, hot}) &= P(\text{sunny} \mid \text{No}) * P(\text{hot} \mid \text{No}) * P(\text{No}) \\ &= \frac{2}{9} * \frac{2}{5} * \frac{5}{14} \\ &= 0.08571 \end{aligned}$$

Now, we have to normalize them to 1

$$\begin{aligned} P(\text{yes}) &= \frac{0.031}{0.031 + 0.08571} \\ &= 0.27 \end{aligned}$$

$$\begin{aligned} P(\text{No}) &= 1 - 0.27 \\ &= 0.73 \end{aligned}$$

So, output is No.

2 Sample Z-test for the difference between Means

A guidance counselor claims that the student at college prep program has higher score than general program. 45 student randomly selected from college prep population with mean score of 24.1 and std of 4.6.

42. Student randomly selected from general program with mean score of 20.2 and std of 4.6.

With 5% significance, can you support the claim

① Hypothesis:

$$\begin{aligned} H_0 &\Rightarrow \mu_1 = \mu_2 & \text{if } \mu_1 = \text{stu in college prep} \\ H_1 &\Rightarrow \mu_1 > \mu_2 & \text{if } \mu_2 = \text{stu in general program} \end{aligned}$$

② Z-test and p-value:

$$\begin{aligned} z &= \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{n}}} \\ &= \frac{(24.1 - 20.2) - (0)}{\sqrt{\frac{4.6^2}{45} + \frac{5.2^2}{42}}} \end{aligned}$$

according to null hypothesis
diff is 0

$$z = 3.695$$

$$p\text{-value} = 1E^{-4} = 0.0001$$

③ ⇒

Decision Rule:

Compare p value to α

$$p\text{ value} = 0.0001 < 0.05$$

Reject the null hypothesis.

It is greater than 0.05

$$S_{\bar{x}} = \text{std. off}$$

$$S_{\bar{x}} < 0.1 < 0.4$$

$$\frac{(x_1 - \bar{x}) + (x_2 - \bar{x})}{2}$$

Z-test for proportion, two samples

64 out of 200 people taking medication report symptoms of anxiety.

Of the people receiving placebo, 92 out of 200 report symptoms of anxiety. Is medication working differently than the placebo. $\alpha = 0.05$

① Hypothesis:

$$H_0: p_1 = p_2$$

$$H_1: p_1 \neq p_2$$

② State Alpha:

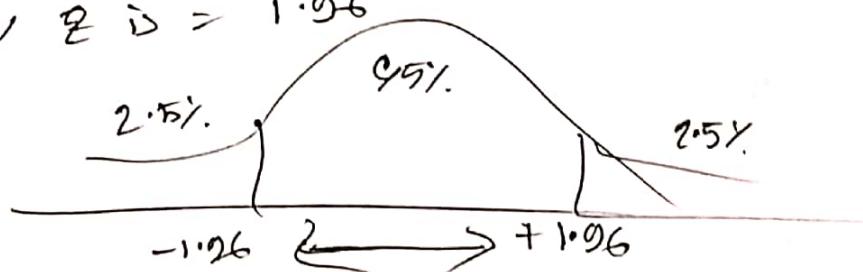
$$\alpha = 0.05$$

③ State the decision:

$$-\alpha = -0.05 = 0.975$$

$$1 - \frac{\alpha}{2} = 1 - 0.025 = 0.975$$

$$60, z \text{ is } = 1.96$$



④ Calculate the test stat,

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$n_1 = 200$$

$$n_2 = 200$$

$$\hat{p}_1 = \frac{64}{200} = 0.32$$

$$\hat{p}_2 = \frac{92}{200} = 0.46$$

$$\text{Z-value} = \frac{0.32 - 0.96}{\sqrt{0.30(1-0.30)}} = \frac{-0.19}{\sqrt{\frac{1}{200} + \frac{1}{200}}} = \frac{-0.19}{\sqrt{0.005}} = \frac{-0.19}{0.0707} = -2.69$$

$\bar{x} = \frac{x_1 + x_2}{2}$

③ State the result:

if Z is less than -1.96 do not reject null hypothesis

$$Z = -2.69$$

\Rightarrow Reject the null hypothesis.

long A short S Q
0.0 < 0.05

$$-2.69 < -1.96$$

Reject the null hypothesis

$$0.0 < 0.05$$

Take the next step

$$\frac{6.9}{8} = 0.8625$$

So the P-value is 0.8625

A/B Testing

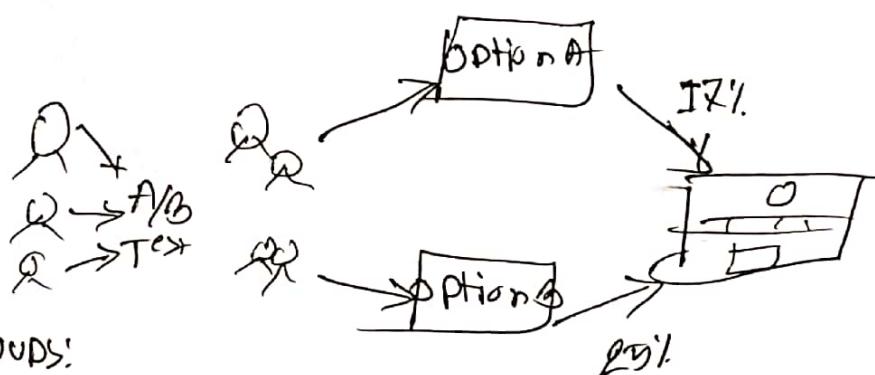
(2)

A

* I want to increase the sell of the product.
Divide the products to A and B. A is unchanged and B has significant changes in packaging.

On the basis of the response from customer groups who used A and B respectively, which one performs better?

Two groups:



① Control

↓
Conversion rate
16%

② Experiment/test

↓
Conversion rate
15%

→ Can we conclude that experiment group is doing better?

No, → for that we have to find that if ~~they are~~ the difference is statistically significant:

through two sample z-test, proportion test

Type 1 error: Reject null hypothesis while it is true. Accept B while it is not performing better than A.

Type 2 error: Accept the null hyp. while it is false. Concluding B is ~~not~~ doing good while it is doing good.