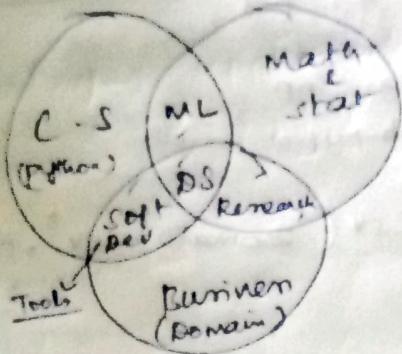


# Statistics → Numerical Science



Data → collection of facts

## statistics

↳ states that "It is a discipline that concerns the collection, org, analysis & interpretation & presentation of data

collect, clean, analysis, Interpret  
(through insights)  
vizual or mathematical  
(patterns expn)

Data is facts  
Data

What is statistics

↓  
Data

↓  
facts (raw facts), especially numerical facts,  
Collected together for reference or information

[Information]

knowledge communicated concerning some particular fact

- Categorical → Numerical → Statistics  
↓  
(other than numerical)

Q      A  
Chi square → Categorical → Categorical  
↓  
Can convert see relationship b/w categorical

Eg: Features (Q)

X

Answer  
Y

# Statistics in real world

Industry → supplement → product  
(for body building) 1 kg  
→ mass, weight, calorie, price

→ daily producer  
1 L products

Industry → pack → 5 biscuits → 1 L products  
other name of sample

① Population  
↓  
1 L product of Biscuits  
Other name of population  
parameter

② Sample →

Random selected Products 1 L  
(i.e., 1000/1 L products)

$$wt = 100 \text{ mg} \rightarrow +12 \text{ mg}$$

$$S.D = \text{Distance} = \sqrt{\text{Var}}$$

mean      data  
SD      median      SD  
↓  
Range  
↓  
spread of data

Central Tendency → mean, median, mode

S.D, Var → Measure of Dispersion (Spread)

1 L → Sample → weight → 95-105 sd.  
(1000 products randomly selected)

Population → Whole set

→ We can't make analysis easily

Sample → Subset of population (loose in groups when compare to the population)

→ Easy to analysis

hence it is a small group

## strategy (sample)

- Sample ~~never~~<sup>is</sup> always resemble the population approximately.

Why sample?

: population → Time, money & hard to do.  
waste

Sample → Easy to analyse

How to take sample:-

i) Random sampling

ii) Intervals

- At frequent intervals

iii) Population

Sample Data should be similar to population.

... few products → Adopt - for Indians  
but not for Americans

Should also consider about climate  
conditions, day/night time,  
Season

Types      Inferential → Resembles, confidence Interval  
                 Descriptive → Describes above the data

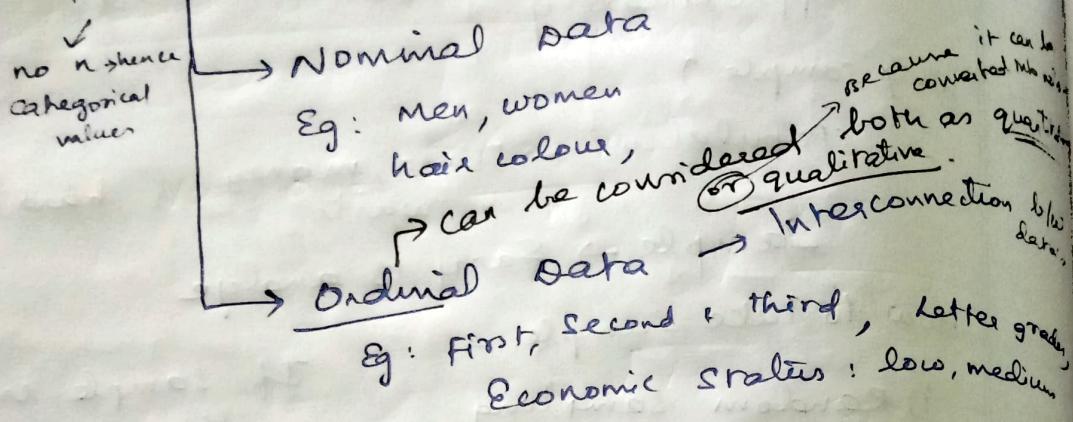
Analysis      Interpretation  
                 presentation

Types of data      i) Quantitative data → numbers  
                          ii) Qualitative data

↓  
Categoricals.

Quantitative      discrete → whole nos  
                      continuous → decimal nos  
                      ↓  
numerical          fractions.  
                      fluctuates

## Qualitative → Categorical values



- Nominal data & ordinal data has to be changed into number. (numerical form)

## A Different types of Analytics: → science of logical analysis

Tells about future  
or predicting a subset division of whole into small components

- i) Descriptive → describing above the problem
- ii) Diagnostic → diagnosing, based on descriptive
- iii) Predictive → prediction
- iv) Prescriptive

## Variable & Random Variable

↳ changes randomly (because of uncertainty)

Measured or counted, they can be weight, height, ...

### i) Numerical variable

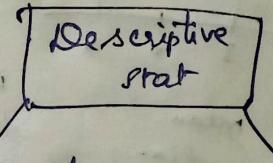
### ii) Categorical variable

Discrete → number  
cont.

## Descriptive

describing - avg, median, sd.

- i) Analysis data - helps to describe & summarize
- ii) Organize, analyse & present data in a meaningful way
- iii) Used to describe a situation
- iv) Explain already known data & limited to a sample or population



Measures of central tendency (mean, median, mode) (with the help of Pandas)

- results shown in charts, graphs.

Data should be in Numerical

Central tendency → @ centre  
(mean, median)

↓  
avg

Sample, mean =  $\bar{x}$   
population, ' =  $\mu$

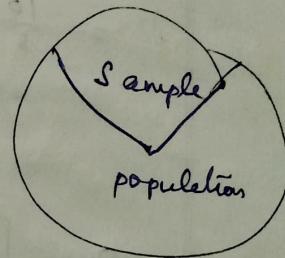
data = 1, 2, 3, 4, 5

$$\text{Mean} = \text{Avg.} = \frac{1+2+3+4+5}{5} = 3$$

median = Middle value

$$\text{even} = \frac{\text{central two terms}}{2}$$

- (researcher's view)
- ## Inferential
- (conclusion based on sample)
- i) Used @ feature selection
  - ii) Analysis of random sample of data & make inference about the population.
  - iii) Compare, test & predict data
  - iv) Used to explain the chance of occurrence of an event.
  - v) Attempts to reach the conclusion



Measures of dispersion  
Variability  
(Variance)  
Range,  
std. deviat.

- Estimation of parameters

- results are shown with prob. scores

Inference abt population with the help of sample

Mode = frequency = repeated values

## Standard Deviation

Measure of Dispersion → spread

$$\text{Var}, \sigma^2 = \frac{\sum (x_i - \mu)^2}{n}$$

$$\text{std. deviation}, \sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{n}}$$

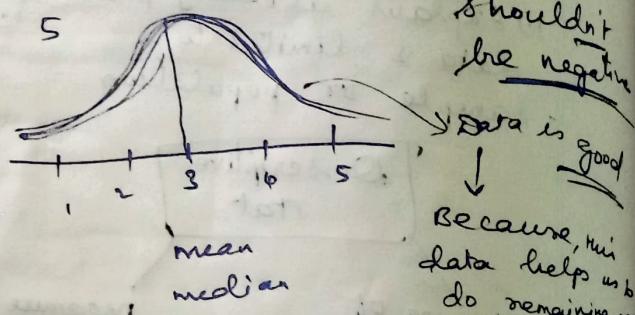
$$\text{std. deviation} = \sqrt{\text{Var}}$$

$$\text{data} = [1, 2, 3, 4, 5]$$

Var =

$$\frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{\text{mean value}} \xrightarrow{2} \text{B/c cause the distance is shouldn't be negative}$$

$$= \frac{10}{5} = 2$$



$$\text{SD} = \sqrt{\frac{4+1+0+1+4}{5}}$$

Real world data  $\Rightarrow$  falls on distributions

Each data follows a pattern

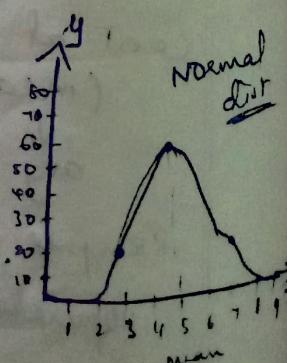
## Normal Distributions

- Data analysis

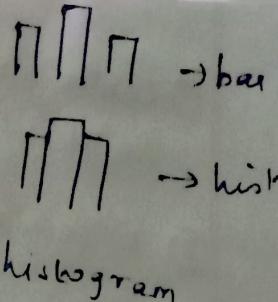
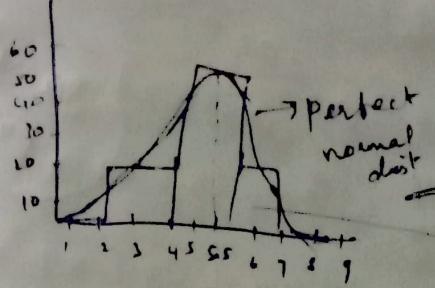
Eg

i) Height of the class

	Small	Normal	Big
Count	20	60	20
Weight	3	5	7



## Histogram



Max<sup>m</sup> pattern in the world follows normal dist.

SD

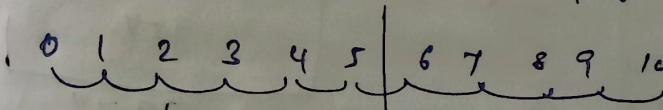
sample

A statistic that measures the dispersion of a dataset → relative to its mean

$$\frac{10}{2} = 5 \text{ math}$$

$$\frac{10}{2} = 5.5 \text{ stat (Numerical Science)}$$

5.5 - 5.6 ⇒ probably blw → physicist

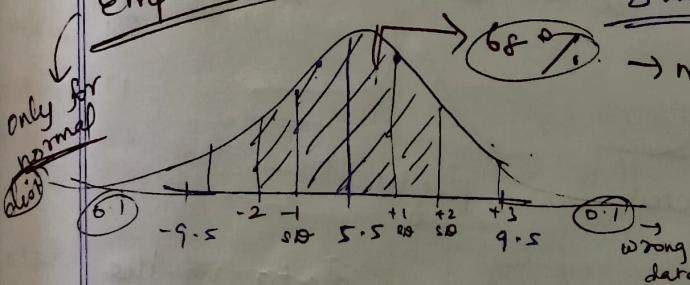


Empirical formula

5.5 → Center

Stat

→ max<sup>m</sup> data holds



Use of studying Central Tendency & Measure of Dispersion

Central Tendency

Mean, median & mode)

avg

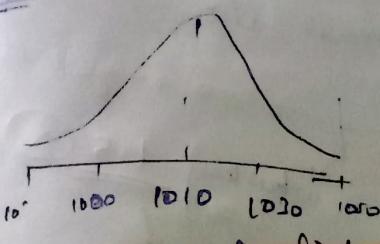
All will be in one straight line

@ centre (O.R. to dispersion)

Empirical formula  
Normal Distribution

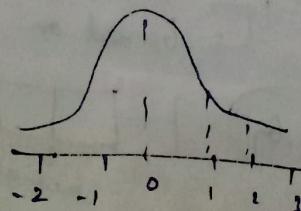
Measure of Dispersion

Variance, SD, Range



normal dist

(No range for its values)



Std. normal dist

mean = 0  
std. deviation = 1

F test → Fisher Test

### Statistics

#### Descriptive stat

- Central Tendency

Mean

median

mode

Measures of Dispersion

- S.D

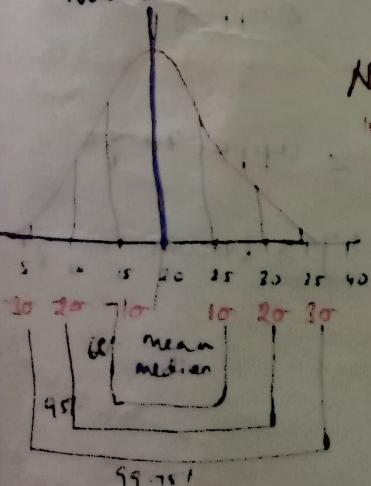
- Variance

- Range

#### Distribution

mean = 20 std = 2

Normal Dist



Normal }  
Dist } Mean = median = mode

⇒ To find outliers ( $> 3\sigma$ )  
⇒ Handle null value  
⇒ Boxplot  
Outliers → Outlier

Boxplot

68% of data fall  
under 1 $\sigma$  std.  
remain

Dot plot → histogram + ESD  
(find = ESD)

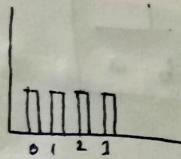
## Types of Distribution



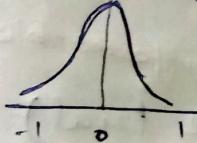
### Discrete

- Uniform distribution
- Binomial → multiple trial
- Bernoulli → only 2 possible outcomes  
In a single trial
- Poisson distribution  
In a controlled environment

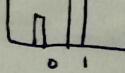
### Uniform.



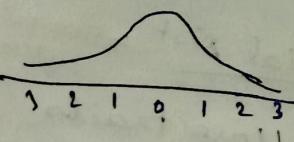
### Std normal dist



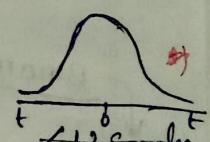
### Binomial Bernoulli:



### normal dist.



### T-dist



Student-t distribution → degree

of freedom → less sample  
(0 to +)

### Binomial

Fair coin → Tail & Head

Bernoulli's distribution → 2 possible outcomes  
only two outcomes

success, Failure  
→ single trial

prob

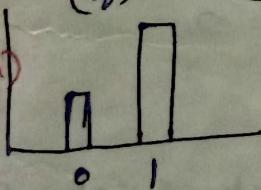
probability  
need not  
to be equal

$$\text{Success} = P = 0.5$$

$$\text{Failure} = Q = 1 - P = 1 - 0.5 = 0.5$$

mean = P

X → random variable



$$\text{mean} \rightarrow E(X) = \sum_{x=1}^2 x \cdot P(x)$$

$$\text{Variance} = P(1-P)$$

$$= PQ$$

$$SE = \sqrt{PQ}$$

$$\text{pmf}$$

Because outcome  
is a discrete

$$P(x) = \begin{cases} 1-P, & x=0 \\ P, & x=1 \end{cases}$$

Success ≠ Failure

prob(failure)

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

$$P, x=1 \rightarrow \text{prob(success)}$$

$$\text{PMF} = \begin{cases} Q = 1 - P, & \text{if } x=0 \\ P, & \text{if } x=1 \end{cases}$$

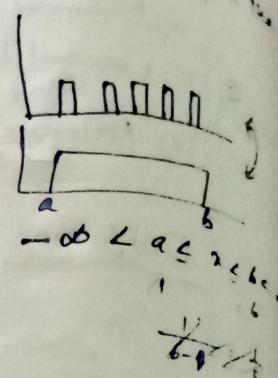
## Uniform Dist

- more possible outcomes  
Eg. Dice (1 to 6)

possible outcomes depends upon material

- same probability

$$f(x) = \frac{1}{b-a}$$



Distribution  
Analysis  
Vidya

expected value = original value

$$(30-15) = \frac{1}{40-10} = \frac{15}{30} = \frac{1}{2} = 0.5$$

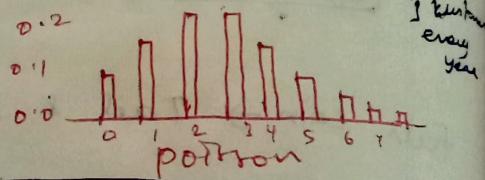
## Binomial Dist $\rightarrow$ Two potential outcomes per trial $\rightarrow$ Difference b/w Bernoulli & Binomial

- More trials (multiple trial)
- Each trial is independent
- Trial (only two possible outcomes)

$P(\text{Success})$  is same across all trials

## Poisson Dist

- discrete distribution
- Describes the no. of events occurring in a fixed time interval or region of space



- Requires only one parameter,  $\lambda$
- Bounded by 0 &  $\infty$

## Assumptions

- The rate at which events occur is constant
- Occurrence of one event does not affect the occurrence of a subsequent event (i.e., events are independent)

Q3(b) coin - 1 time toss  $\rightarrow$  Bernoulli ✓  
" - more time toss  $\rightarrow$  Binomial ✓

### Poisson

- parameter,  $\lambda$  (Expected rate of occurrence)
- Prob. of given no. of events in a fixed time interval (within a given time interval)

Eg: no. of calls attend in 1 hour at call center.

### Rules

① 1 success  $\rightarrow$  2 success (not influence)  
1 hr (30 calls)      35 calls  
2nd hr

② Small Time  $\rightarrow$  Success rate  
long time  $\rightarrow$  Success sum

small time (1 hr) =	Raja	Raja
25	30 ✓	
10	20 ~	
— 35	20	

$P(S)$  over a short interval is equal to the  $p(S)$  over a long time interval because attended more calls.

③ Interval becomes smaller because the interval approaches zero.

$\rightarrow$  prob. Density fn

PDF  $\rightarrow$  Continuous  $\rightarrow$  parametric (normal dist)

PMF  $\rightarrow$  Discrete  $\rightarrow$  non-parametric  
 $\downarrow$  discrete value  $\hookrightarrow$  chi  $\chi^2$  dist

prob. Mass fn

student t test

Anova  $\rightarrow$  Z test + F test  
(F test)      1 on 1      1 on 1

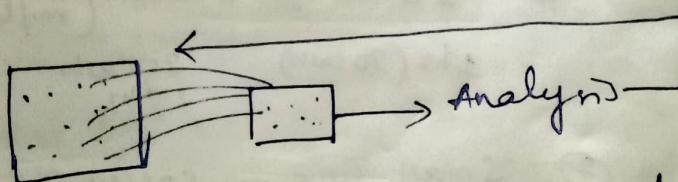
$\downarrow$   
Compares more than 1

# Inferential statistics

To draw inference about the population from sample data.

Inference → Conclusion reached on the basis of evidence and reasoning.

To draw inference about the population data from sample data.



describing data → descriptive stat

Taking sample from population & concluding → Inferential stat.

population → sample → Analysis

Inferences about population → Result

If it should resembles the population.

$t$  test,  $f$  test, chisq

1<sup>st</sup> time result check → Null hypothesis

2<sup>nd</sup> time " " → alternate hypothesis  
(Researcher hypothesis)

i) Reject null hypothesis

ii) Failed to reject null hypothesis

Assumption → No differences → Null hypothesis

Difference → (Accept / Failed to reject null hypothesis)

# Confusion Matrix

		Actual (real)	
		+	-
predicted	By model -	T	F
		T	TF
		FT	FF

Type-I Error (FP) → True Positive  
 Type-II Error (FN) → True Negative  
 (Signon +ve)  
 (Signon -ve)  
 False Positive  
 False Negative

P		
P		
I		

Actual | predict → (should consider this one).

$$T \quad T = TT$$

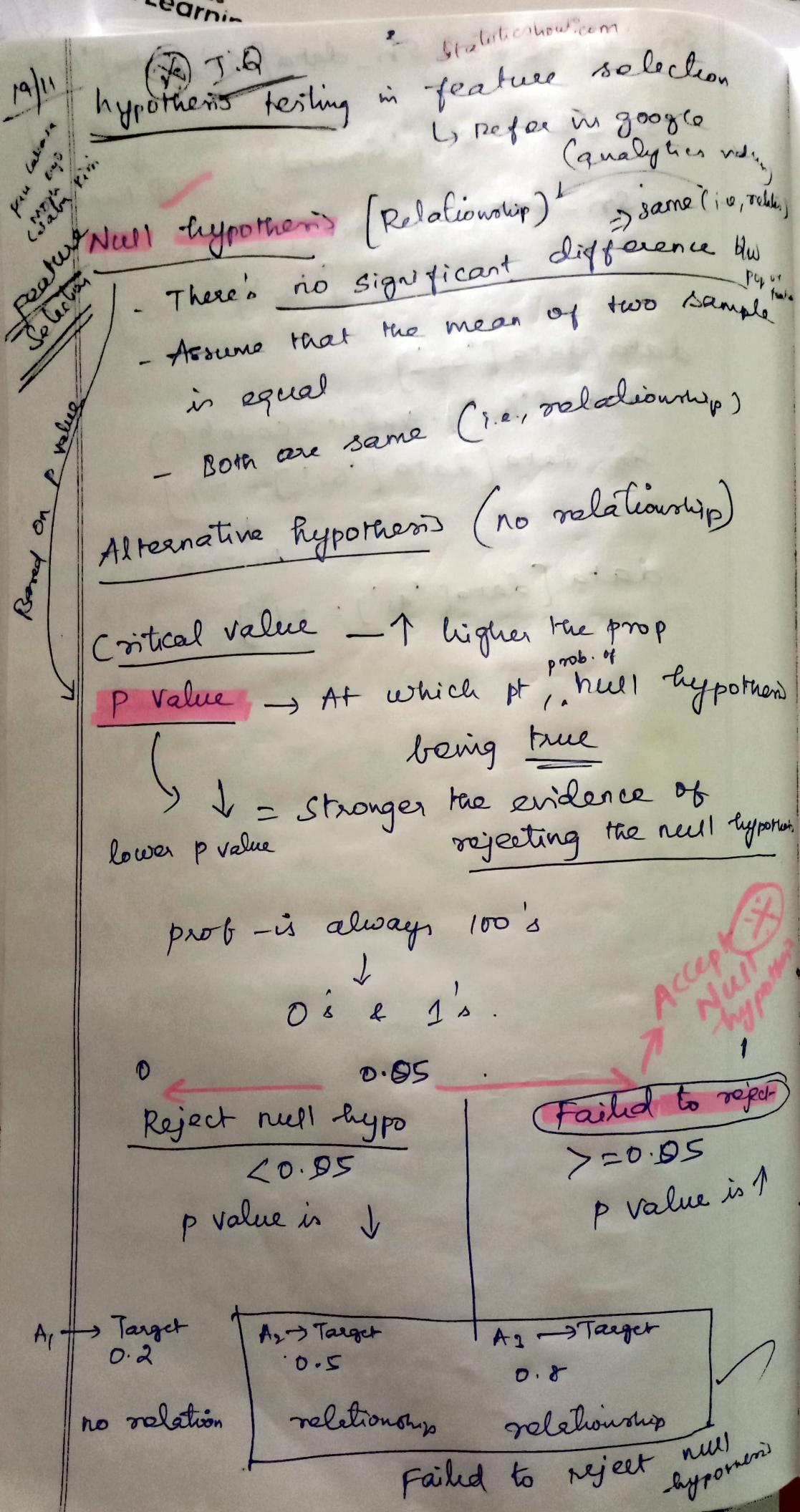
$$F \quad F = FF$$

$$F \quad T = \text{False T} \rightarrow \text{False +ve}$$

(negative 2) Example

$$T \quad F = \text{False F} \rightarrow \text{False}$$

negative  
(positive 2) Example



7. 96
1. One-on-one prediction → for public
  2. Batch prediction → Industries, hospitals.
  3. "Run linearly once"
  4. Streaming → Stock Market

minimizing  
for  
scaling  
the  
data

### curve of dimensionality

multi-collinearity → features ~~Corr~~  
multi-collinearity ~~Corr~~ its effect

$$\text{Confidence Interval} = 100 - p \text{ value}$$

Anova → can compare more than ~~two~~ variables

### Hypothesis Testing

#### Statistics for data science

→ p value, N.K.Bi → study abt hypothesis,  
Type-I & Type-II

#### → Biostatistics

##### Type-I error :

(False-positive) → actual false predicted True

→ occurs if an investigator rejects a null hypothesis that is actually true in the population.

##### Type-II Error (False-Negative)

- occurs if the investigator fails to reject a null hypothesis that is actually false in the population

Inferential  $\rightarrow$  hypothesis testing

feature  $\rightarrow$  Target  $\Rightarrow$  should have relationship  
B/w two features  $\Rightarrow$  NO relationship  
should be there.

p value  $= 0.05 \Rightarrow$  may get varied

$\chi^2$  test  $\Rightarrow$  Converting normal dist into  
std. normal dist.

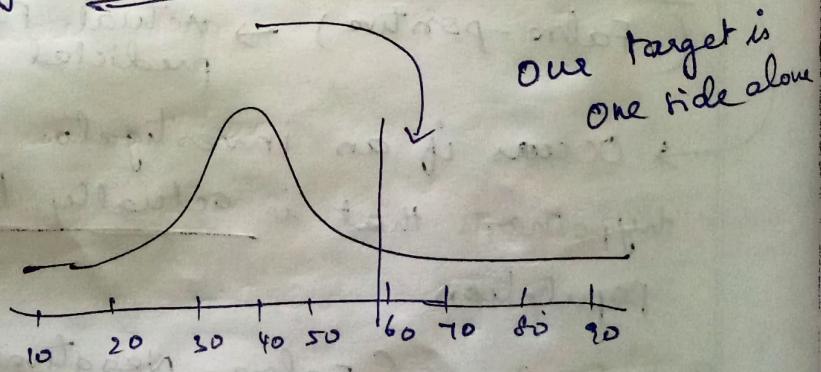
Level of Significance -  $\alpha$  (alpha)

$\alpha = 5\%$  risk  $\Rightarrow$  okay to take  
5% risk  
Very for  
medical data  
 $(\alpha=0.01)$   
 $= 0.05$  error are allowed.

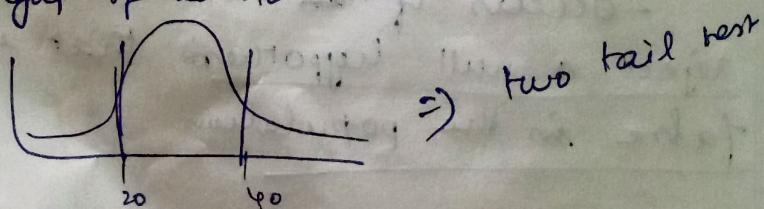
Critical Value (C)  $= 100 - \alpha$

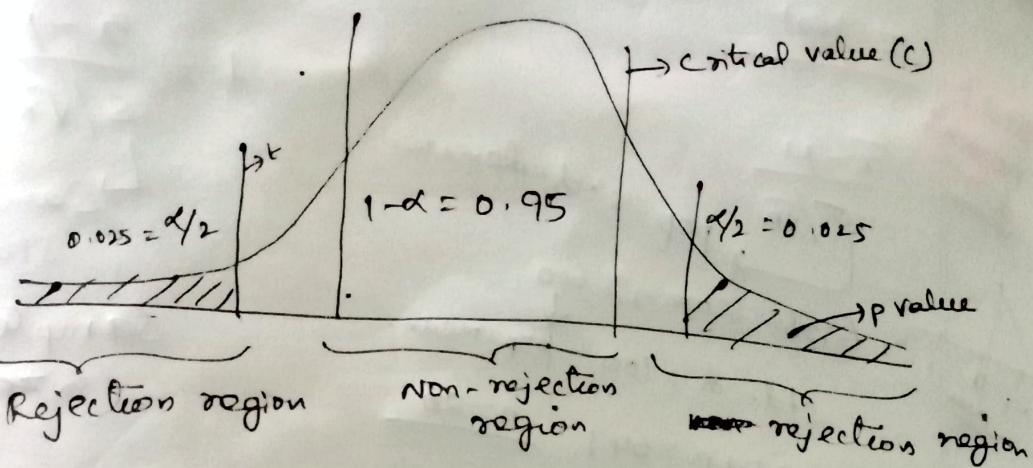
p-value  $\rightarrow$   
single-tail test  $\rightarrow$  focusing on single side  
two-tail test

i) ppl got 60+ marks  $\rightarrow$  score high  $\Rightarrow$  null hypothesis



ii) age gap of 20-40 smokers low





$t > c \rightarrow$  reject  $H_0$

$t \leq c \rightarrow$  failed to reject  $H_0$

$P > \alpha$	$\rightarrow$ fail to reject $H_0$
$P \leq \alpha$	$\rightarrow$ reject $H_0$

### Hypothesis

1) Define null, Alternative

2)  $\alpha$

$$\begin{aligned} \downarrow \\ \text{allowed error} \end{aligned}, CI = 1 - \alpha \quad \begin{aligned} \xrightarrow{\% \text{ proportion or mean}} \\ \xrightarrow{\text{one tail or two tail?}} \end{aligned}$$

$$\begin{aligned} &= 1 - 0.05 \\ &= 0.95 \end{aligned}$$

3) score

- critical value  $\rightarrow$  which test have to choose?  
 $Z, t, f, \chi^2, \dots$

Eg:  $18 \pm 2 \Rightarrow 16, 20$

$\Leftrightarrow$  b/w null  $H_0$  is true

A) p value

- prob of being  $H_0$  is true

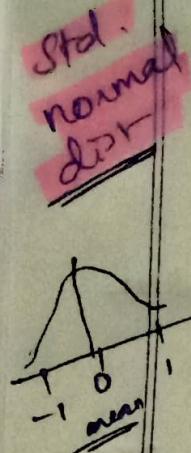
	$H_0$	$H_a$
$H_0$	No error	Type I error
$H_a$	Type II error	No error

## T test

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

Sample  
std. deviation  
Sample size }  $\geq$   
 $< 30$

If it is greater  
than 30, Use  
Z-test



## Z test

$$\frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

pop. std. deviation  
Gird's score  $\Rightarrow 600$   
std. d = 100      ↓  
mean = 641      pop. mean.  
 $n = 20$

$$Z \text{ score} \Rightarrow \frac{641 - 600}{100 / \sqrt{20}}$$

$$= 1.8336.$$

In z table,

$$\underline{p \text{ value} = 0.033}$$

we reject  $H_0$

( $\because p \text{ value} < 0.033$ )

Hence the given values  
are false.