

another level of stats

Statistical methods

Descriptive

Univariate
↳ shape
↳ center
↳ spread

Bivariate
↳ Correlation
↳ Regression

Multivariate
↳ Heatmap
↳ Multiple Linear Regr

Inferential stats

↳ applied to mean
(Z-test, t-test)

↳ applied to Variance
(F-test, ANOVA)

How many way to sum
SCT

Ratio \rightarrow eg
• Male female ratio meas
• Cloths 2:1
• Temp $\frac{21}{10} \rightarrow$ inside 30° out 60° \rightarrow 1:1

\rightarrow Mukhya ji Mean
In you in USA
Caste Indian (mean)
all work around the
mukhya

3 obj \rightarrow mean = 3 1, 2, 3, 4, 5 - obj

mean = 3, 3, 3, 3, 3 - obj

\rightarrow Describe of data at 2nd Both objs same? Now you stop

\rightarrow So do you really think all says about data \rightarrow No

\leftarrow that's why the another topics comes

② Measure of Spread or dispersion

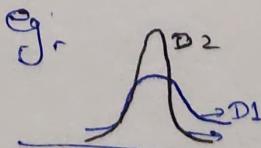
eg obj: O1 { 1 2 3 4 5 }
D2 { 3 times }

So you understand
why spread use
is
bcz alone CT
not enough info of data

3

Then why go for Measure ~~of~~ Symm

③ Measure of Symmetry



Spread of D1, D2 same for both
but the shape is difference so
that's why learn (Skewness kurtosis)
for high / short tail distribution

۹

mean battlefield area %

Model

1,2,2,3,3,45

→ called bi-model
means (2) models
2, & 3

② 2, 2, 3, 3, 4, 4, 5, 6, 7

2, 3, 4 called Multi Model

—m

\Rightarrow Arithmetic mean $\rightarrow \sum_{j=1}^n \frac{x_j}{n}$

Q) **Geometric mean** → multiplies the no together and does a square root case so on
 Ratio same or d. how many number are there.
 Eg: 2 & 10 → GM $\Rightarrow \sqrt{2 \times 10} = (2 \times 10)^{\frac{1}{2}} = \sqrt{20} = \text{G.M}$

$$\text{eg } \textcircled{1} \quad 2\sqrt[3]{25} \xrightarrow[\substack{\text{no} \\ \text{method}}]{\substack{3 \\ \text{no}}} \sqrt[3]{2 \times 25} = \sqrt[3]{(2 \times 25)^2}$$

Because
eg If you go for buy a Camera
C1 way: Prob zoom 20x
11 25D

- GM is used when we want to compare two different properties

$$C_1 = \sqrt{200 \times 8} = 210 \rightarrow \text{BIM high so Select } \\ \sqrt{250 \times 6} = 38$$

③ Harmonic Means

$$H \cdot M = \frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \dots + \frac{1}{n}$$

$$\text{eg: } 2, 4, 5, 100 \Rightarrow \frac{4}{\frac{1}{2} + \frac{1}{4} + \frac{1}{5} + \frac{1}{100}} = \frac{4}{0.5 + 0.25 + 0.2 + 0.01} = \frac{4}{1}$$

↓
is the reciprocal of avg of all reciprocals.

$$\text{Mean} \downarrow \\ \text{avg. of seci} = \frac{\frac{1}{1} + \frac{1}{2} + \frac{1}{3}}{3} \rightarrow \text{recip recip} = \frac{3}{\frac{1}{1} + \frac{1}{2} + \frac{1}{3}}$$

Use-in

① Confusion matrix \Rightarrow F. score (F-score is nothing but H.M of recall & precision)

Q2) we travelled 10km at 60km/hr, then another 10km at 20km/hr. What is avg Speed?

$$4 \sqrt{\frac{2}{Y_{60} + Y_{20}}}$$

ed? may be as real in data theory/CI

④ Frequency Means (f_n)

e.g. score frequency

		\rightarrow 2 times
1	2	
2	5	
3	3	
4	6	
5	5	
6	2	then what

$$\Rightarrow \boxed{\frac{\sum f x}{\sum f}}$$

then what is Mean of $\frac{29}{29}$ Score (indicated add time taken)

* Weighted Mean :

if calculator mean = $\frac{1+2+3+4}{4}$ \rightarrow what actually happen here \rightarrow avg here but actual

So we can say

All the values have equal weights

1 has 1 freq
2 " 1
3 " 1

1 2 3 4

El Number
Weight importance
deter.

here equal wt
so $\frac{1}{4}$ equal
all

So we can write this in this form

$$1 \times \frac{1}{4} + 2 \times \frac{1}{4} + 3 \times \frac{1}{4} + 4 \times \frac{1}{4}$$

Now the Question is here

e.g.: 1, 2, 3, 4 \Rightarrow calculate the mean $= 0.1 \times 1 + 2 \times 0.1 + 3 \times 0.7 + 4 \times 0.1$
weight $\rightarrow 0.1, 0.1, 0.7, 0.1$ How $= 2.8$ means

Use Case / Application

- Weight mean can be used where something are important and as compared to others.

e.g.: $[1] [0.1] [1] [0.1]$ \rightarrow what is sum of weights $\rightarrow 1$ if one the not need to divide we remove

eg(1) Camera example :- Before
Real world example
wants to rank by \rightarrow Image Quality : 50%,
Battery life : 30%,
Zoom range : 20%.

Cam 1	Reviews	Cam 2
Sony \rightarrow IQ: 8	BL: 6	IQ: BL: ZR
ZR: 7	9 4	6

do here W-Mean \rightarrow Which camera should be Reference
 \rightarrow So what you do Mean or AM, GM, WM
 \rightarrow Sony $\Rightarrow 0.5 \times 8 + 0.3 \times 6 + 0.2 \times 7$ | Cam 2 $\rightarrow 0.5 \times 9 + 0.3 \times 4 + 0.2 \times 6$
 $\Rightarrow 7.2$ | $\Rightarrow 6.9$
 prefer

eg(2) Dimesh had lunch 7 times a week but some weeks he get only 1, 2 or 5 lunches. So on 2 weeks \rightarrow only 1 lunch for whole week

$$\begin{array}{l} 14 \\ 8 \\ 32 \end{array} - \begin{array}{l} 2 \\ 5 \\ 7 \end{array}$$

Ask \rightarrow then what is the mean no of lunch that Dimesh have every week ??

Use simple mean \times not proper o/p

weight is more than 1
 $\rightarrow 50 \rightarrow 2+14+8+32 \rightarrow 56$
 \rightarrow so divide

$$\frac{2 \times 1 + 14 \times 2 + 5 \times 8 + 32 \times 7}{56}$$

$$\Rightarrow \frac{56}{56}$$

Recap \rightarrow Summary
descriptive \rightarrow 1, 2, 3;
 SL - 1, 2, 3, 4, 5; 3, 3, 3, 3;
 MGT 3, 3 \rightarrow 3, 3, 3
 MCT \rightarrow not alone true
 \rightarrow Spread of data
 Variance, Range, Percentile/Box Plot
 Quartile (IQR, Box Plot) \rightarrow Var.

⇒ Dispersion

① Range → difference b/w max^m & min value

e.g.: 1, 2, 3, 4, 5

$$\text{range} = 5 - 1 = 4 \Rightarrow \text{It is known from data so what now}$$

Note:-

+ outliers affect the range

e.g.: 1, 2, 3, 4, 1000

$$\left. \begin{array}{l} \text{range} = 1000 - 1 \\ \Rightarrow 999 \end{array} \right\} \begin{array}{l} \text{here calc is correct} \\ \text{but ambiguous interpretation} \end{array}$$

② Percentage & Percentile:

e.g.: 1, 2, 3, 4, 5

$$\text{odd \% } \frac{3}{5} \times 100 = 60\%$$

↳ It is a value below which a certain % of observation lie.

e.g.: data = 1, 2, 3, 4, 4, 6, 7, 8, 10

↳ def. described
Show percentiles

③ What is the percentile rank of 3?

$$\text{ptile rank} = \frac{\text{No of Value below that number}}{\text{Total number}} \times 100$$

$$= \frac{2}{10} \times 100 = 20^{\text{th}} \text{ percentile}$$

↳ What value exist at 75th percentile

↳ Means 20%. non-eq to or below 3
20% तक समान नहीं हो सकता

$$\text{Value} = \frac{\text{percentile}}{100} \times (n+1) \rightarrow \frac{75}{100} \times (10+1)$$

$$= \frac{75}{100} \times 11 = 8.25^{\text{th}} \text{ no.}$$

The 8th no that is 7 is my 75th percentile

↳ So take avg of 8 & 9
8th & 9th no

8.5 → 8
8.7 → 9

③ Quartile ; Quartiles are values that divides a list of no into Quarters (4 equal part)

Solve → Pages Outliers

→ e.g.: 6, 8, 5, 5, 7, 3, 9

↳ order the no :- 3, 5, 5, 6, 7, 8, 9

3rd | 2nd | 3rd | 4th

Q₁ - 5
Q₂ - 6
Q₃ - 8

↳ if bigger no. of data sets

4 Quarters / 3 Quartiles

e.g. Jan feb ... dec 1, 1, 2, 2, 3, 3, 4

e.g. 1, 1, 2, 2, 3, 3, 4

↳ what is Quartile 3?

① Order

② total count = 11

$$\rightarrow \text{if count is odd} \rightarrow Q_1 = \frac{n+1}{4}^{\text{th}} = \left(\frac{n+1}{4} \right)^{\text{th}} = 3^{\text{rd}} \text{ obs} \rightarrow$$

$$\rightarrow Q_2 = \frac{(n+1)}{2}^{\text{th}} = \left(\frac{n+1}{2} \right)^{\text{th}} = \left(\frac{11}{2} \right)^{\text{th}} = 6^{\text{th}} \text{ obs} \Rightarrow 2$$

$$\rightarrow Q_3 = \frac{3(n+1)}{4}^{\text{th}} = \frac{3 \times 11}{4}^{\text{th}} = 3 \times 3 = 9^{\text{th}} \Rightarrow$$

= (if even)

→ Q₁ = $\frac{n}{4}$ th observation

→ Q₂ = $\frac{3n}{4}$ th

→ Q₃ = $\frac{3(n+1)}{4}$ th + $\left(\frac{n+1}{2} + 1 \right)^{\text{th}}$

$$\text{e.g.: } 1, 2, 3, 3, 4, 4, 4 \rightarrow Q_1 = \frac{6}{4} = \frac{3}{2} = 1.5 \rightarrow \frac{1+2}{2} = 1.5 \text{ Q}_1$$

$$Q_2 = 4.5 \text{ m}$$

$$\frac{4+4}{2} = 4$$

→ Understand the importance of Quartile:
 ↗ need
 ↗ DS/DA

↳ def. describe) → FIVE Point Summary

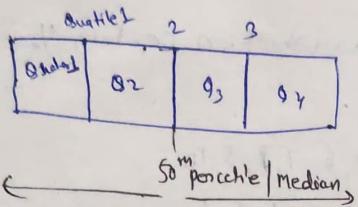
$Q_1 \rightarrow 25^{\text{th}}$ min → 1st Quartile
 $Q_2 \rightarrow 50^{\text{th}}$ " → Median
 $Q_3 \rightarrow 75^{\text{th}}$ " → 3rd Quartile
 $Q_4 \rightarrow 100^{\text{th}}$ max → 4th Quartile

Eg.: Transaction amounts

1500
1200
600
700
2

Let say $Q_1 = 1000$
 $Q_2 = 2000 \rightarrow 50\%$ of the population is the transaction of
 $Q_3 = 3000 \quad 2000 \text{ Rs or below } 3000$

→ Box Plot:



↳ What is Outlier (extreme value → Much much higher, much much lower)

Eg. (-100) 1, 2, 3, 4, 5 (100)

So how do we calculate the Range without affecting the Outliers.

where is possibility to lie outlier in Quartile one → Yes (1) & (4)

Come here

$$\text{Inter Quartile range} = Q_3 - Q_1$$

Eg.: 2, 3, 3, 3, 3, 4, 4, 5, 5, 5, 6, 6, 6, 7, 8, 9, 9

$$Q_1 = 25^{\text{th}} \text{ percentile} = \frac{25}{100} \times 16 = \frac{1}{4} \times 16 \Rightarrow 4^{\text{th}} \text{ row} \Rightarrow Q_1 = 3$$

$$Q_3 = 75^{\text{th}} \text{ percentile} = \frac{75}{100} \times 16 = \frac{3}{4} \times 16 \Rightarrow 12^{\text{th}} \text{ row} \Rightarrow Q_3 = 6$$

$$\text{IQR} = Q_3 - Q_1 = 6 - 3 = 3 \text{ Ans}$$

Useful to find Outliers

In normal distribution $Q_1 - Q_3$ (68% lies)

Points
 → Most data in Q1-Q3

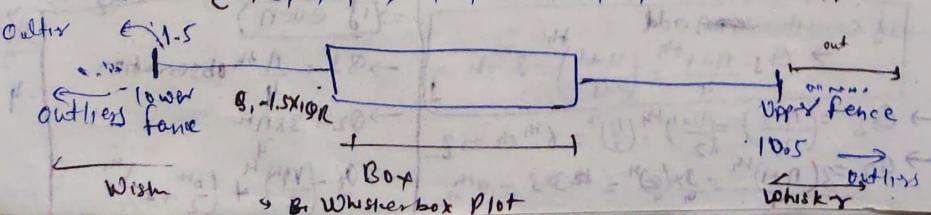
How

$$\text{Lower fence} = Q_1 - 1.5 \times \text{IQR}$$

$$\text{Upper fence} = Q_3 + 1.5 \times \text{IQR}$$

Why = $1.5 \times \text{IQR}$
 any number
 but research conclude

Eg.: {2, 3, 3, 3, 3, 4, 4, 5, 5, 5, 6, 6, 6, 7, 8, 9, 9}

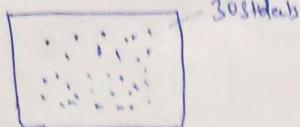


Topics
 Mean Deviation
 & Variance
 ④ Std Dev
 & Measure of Symmetry
 • Skewness
 • Kurtosis

* Covariance & Correlation
 * Probability

→ Measure of Spread → Variance & Std dev

Eg:



In School Auditor/Principals come to school, he don't know about class, and he wants to judge the student's how they perform in class. Compare to others.

• Infrastructure → he give SMCB → 30 Question

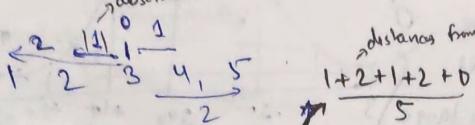
then calculate Avg Mark = 2.1

Data 1, 2, 3, 4, 5 and take the No. of student if the 2.9 & good 1.5 is below

How away each data away from the mean?

Cal → mean = 3

→ away absolute value



$\frac{1+2+1+2+0}{5} = \frac{6}{5} = 1.2 \rightarrow$ each point Avg Away from Mean = 1.2

Q again → Reform Quesn → on an Avg how much away each data point is away from mean Value ??

$$\sum_{i=1}^n |x_i - \bar{x}|$$

mean deviation

Variance ; The average of the Squared differences from the mean.

$$Var_{pop} (\sigma^2) = \frac{\sum (x_i - \bar{x})^2}{N}$$

$$Var_{sample} = \frac{\sum (x_i - \bar{x})^2}{n-1} (S^2)$$

① Calculate Mean

Changes N - n

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

$n - n-1$ (denominator)

→ for each no in data, subtract the mean & no.

→ Square the difference

→ Cal the Avg square of diff

og

$$data = \{1, 2, 3, 3, 4, 4\}$$

$$x_i - \bar{x} = (x_i - \bar{x})^2$$

$$\begin{array}{|c|c|c|} \hline i & x_i & x_i - \bar{x} \\ \hline 1 & 1 & -1.33 \\ \hline 2 & 2 & -0.33 \\ \hline 3 & 3 & 0.33 \\ \hline 3 & 3 & 0.33 \\ \hline 4 & 4 & 1.33 \\ \hline 4 & 4 & 1.33 \\ \hline \end{array}$$

$$\sum (x_i - \bar{x})^2 = 2.83$$

$$Var_{sample} = \frac{6.82}{n-1} = \frac{6.82}{5} = 1.37$$

What is actually it? find what is spread data

① σ ② \bar{x} ③ $\frac{1}{n}$
 what
 → Var ↑ spread ↑
 → Var ↓ spread ↓

* Standard deviation : Standard deviation is a measure of how spread out numbers are.

b) Square root of Variance

$$Std = \sqrt{Var} = \sqrt{1.37} \Rightarrow 1.17$$

$$Replica : SD of pop = \sigma = \sqrt{Var_p}$$

$$SD of sam \Rightarrow S = \sqrt{Var_s}$$

Number seven
when compare
other

Use Case

Math is fun examp

$$\text{standard deviation}^2 = 21704$$

$$\sigma = \sqrt{21704}$$

347

Stand Dev is a standard way of knowing how normal, large extra large is something

$$\text{mean} = 394$$

$$SD = 147$$

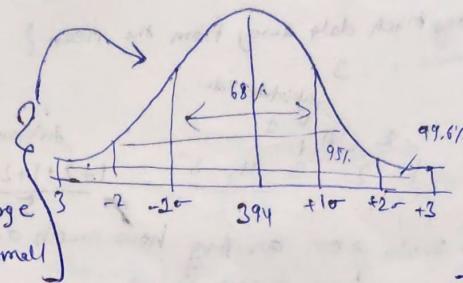
$$\rightarrow 1SD \Rightarrow 394 + 147 = 541$$

$$394 - 147 = 247$$

Example

$$2SD \Rightarrow \begin{cases} 394 + 2 \times 147 \Rightarrow \text{large} \\ 394 - 2 \times 147 \Rightarrow \text{small} \end{cases}$$

$$3SD \Rightarrow \begin{cases} 394 + 3 \times 147 \Rightarrow \text{extra large} \\ 394 - 3 \times 147 \Rightarrow \text{extra small} \end{cases}$$



In Real life
the Variance not use
only show spread
it shows understand

Useful SD comp so
we compare the
data

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

if wrong sample
because we cannot access pop in

why $(n-1)$? Basel correction / Unbiased estimation

We use $n-1$ rather than n is because Sample Variance will be unbiased estimating

→ How much bigger Sample Can be ??

pop → Sample →

pop → subset → Sample → 10 ले हो कि तो Sample become population

$10-1 = 9$ ज्यादा atleast 1 कम हो $\frac{n-1}{n}$

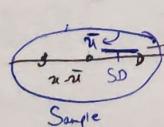
e.g.: $N=10$ अगर Population size 10 हो तो 10 ले हो Sample लेंगी तो Pop & Sample

X

but it is not
correct.

Completely correct.

really
eg! Pop



Step → Pop Data → Pop Mean
Select → Sample → then \bar{x} → why जबकि sample लिया है?

N → यदि यह Pop है तो Sample Variance calculate करना है at Sample 1st

$\text{PopVar} = \frac{(x_i - \bar{x})^2}{N}$ → but you don't have a acc of pop

→ Pop Averages नहीं हैं sample तो M diff
→ at PD जाये हैं SD

→ क्योंकि $(x_i - \bar{x})^2 > (x - \bar{x})^2$ → क्योंकि यह Assumptions के गलती Point के लिए true होता है।

So that's why

$$\frac{(x_i - \bar{x})^2}{N} > \frac{(x - \bar{x})^2}{n-1}$$

$n = \frac{10}{5} = 2$ Ratio रहा और 10 → 8 बढ़े हो तो denominator

Ratio same हो जाएगा।

denominator को

$\frac{8}{4}$

जहां कम होता है वह 2 Constant रहता है जो समान है।

यदि $x_i - \bar{x}$ के Value of को आपना लिये यह Sample का है तो
यदि $x - \bar{x}$ के रखते Ratio same रख दें तो लिए $n-1$ लिये

∴ If PD, SD & ratio of maintain रखना है तो क्या बताएं S.Var में n-2 क्यों लगते हैं P

ताकि $\frac{(x_i - \bar{x})^2}{N} > \frac{(x_i - \bar{x})^2}{n-1}$ True हो

but why then $n-2$ not $n-2, n-3, \dots$

→ Concept here → **Degree of freedom** [total no. - 1] max frequency
maximum no of logical independent Variables.

Eg: 6 5 9 2 28
 $\frac{6+5+9+2+28}{5} = \text{avg of } 5 \text{ no is } 10$
 ↑ ↑ ↑ ↑ ↑
 Any Number what no 9 but here constant

Eg: 100 → 100 → तो यही नहीं किसी No. की At least fix रखना होगा
 ↗ No of possible are Changeable

→ 1 element should not be changeable

So that $n-1$ come into picture

* * * So max^m reduction you make $n-1$ to be called as a sample

• You remove 1 element and only 1 part called as a sample

Measure of Symmetry:-

① Why learn? what is

→ Then how to measure the Symmetry?

Helping ① Skewness

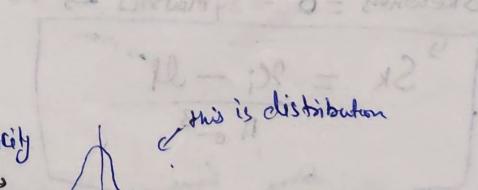
② Kurtosis (shape lit.)

→ Symmetry
 Object syn
 mirror image

① Skewness → dataset's symmetry

If the given data is symmetric e.g.

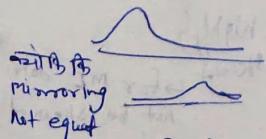
In nature then Skewness = 0.



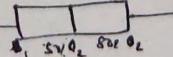
called symmetric ⇒ Skewness
 Mean = Median = Mode

↑↑ why → because skewness is not symmetric.

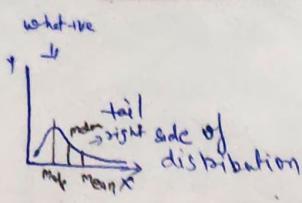
Then how skew data look like
 non-symmetric



→ If in boxplot →



(a) Positive Skewed : (Right Skewed data)



So where where is Mean Median Mode
mean > Median > mode

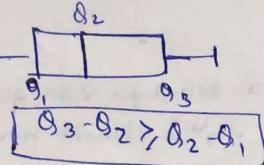
median = physical
mid point

→ What is the reason of Right Skewed?

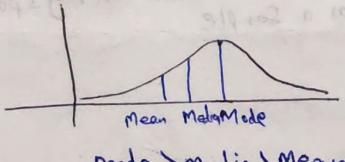
- due to outliers (परस्पर नहीं होते तो यहाँ कोरेक्ट नहीं होता रहता उसका कोई sense नहीं आ रहा।) → e.g.

e.g:

→ Box Plot



(b) Left Skewed data



mode > Median > Mean

Reason → outliers

e.g. ① Death rate before 50 age

② Marks in easy exam

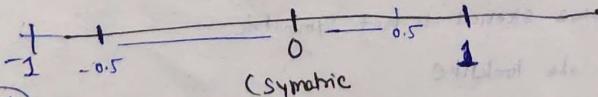
③ Wealth distribution

Q How Skewness = 0 → Symmetric \rightarrow data

$$Sk = \frac{\sum (X_i - \bar{X})^3}{n \sigma^3}$$

de no Sigma

if Skewness



-0.5 - 0 } fairly symmetric

0 - 0.5 }

- (-1) - (-0.5) } moderately skewed

0.5 - 1 } skewed

worry about

< -1 & > 1 } highly skewed

for ML data

but be skewed

+ because of outlier

!!

Question → why Outliers bad

in ML

some

that not give you

correct predictions

→ How to see Skewness / How to know

by ① Visualization

① O, Q Plot

③ Skewness formula (Statistical way)

→ How to use in Industry

• Visualization

→ How to treat Skewness

↳ Transformation

① treat outliers

② Box-Lox transformation / geo Johnson

③ Exponential transf

④ Reciprocal

⑤ log transform

why all we learn stats

↳ 1st moment

$$\text{Mean} = \frac{\sum (x_i)}{N}$$

What common in these

2nd moment

$$\text{Var} = \frac{\sum (x_i - \bar{x})^2}{N}$$

mean is computing taking 0 as reference

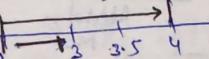
3rd moment

$$\text{Skewness} = \frac{\sum (x_i - \bar{x})^3}{N^{3/2}}$$

4th moment

$$\text{Kurtosis} = \frac{\sum (x_i - \bar{x})^4}{N^2}$$

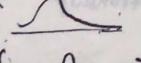
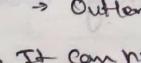
What is actual meaning of mean?
↳ $\frac{\sum (x_i - \bar{x})}{N}$ → On an avg how much a specific value or data point lie away from origin.

↳ if Avg of 2 no. 3 & 4 → 3.5 → 

Question

me say

→ Outlier

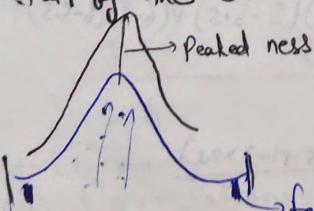
↳ If  → It can not have outliers?  X not necessarily

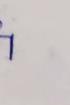
if we can't find outliers in Visualization
then we **Box plot**.

↳ Kurtosis

$$a_4 = \frac{\sum (x_i - \bar{x})^4}{N^2}$$

↳ tail of the distribution / fattness of the tail.



fatness of tail → but the peakedness it depend 

Now Change

• $k > 3$ → Leptokurtic → tail is fat → Many outliers

• $k = 3$ → Meso

• $k < 3$ → Platykurtic & also check is it ok.

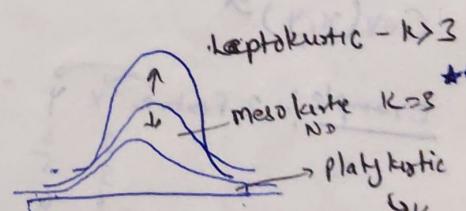
k is reference point

Use case!

• Buy Stock in Bulky

• Mutual funds house

• Sudden Spike in data (crypto) 2022 $k < 3$ we in safe port



Correlation / Covariance!

e.g. ① House price feature \rightarrow house

Coln \hookrightarrow Whatever you are going to predict / Price of a house, it should have some relation with the features.

\hookrightarrow Price of house & Area of house, No of rooms, locality, relationship -

② No of match played \rightarrow & total score

Q How to measure these relationship?

① Covariance

② Correlations

③ which feature we use or not in ML.

Ans: Features which has high correlation value with target variable not going to be used.

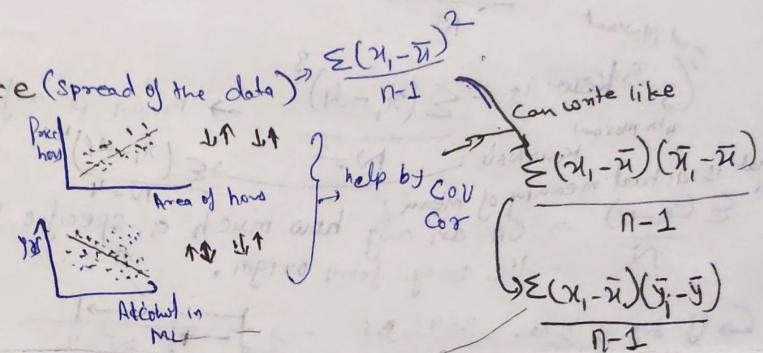
* Covariance.

What is Variance (spread of the data) $\rightarrow \frac{\sum(x_i - \bar{x})^2}{n-1}$

CO = Variance

\downarrow before it $\xrightarrow{\text{Scenario 1}}$ Scenario 1
 \downarrow Scenario 2

To determine relationship b/w 2 variables.



So we can also say Variance is nothing relationship itself.

$$-\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$\text{Cov}(x,y)$

Example: 2 Feature $x \quad y$
 $\begin{array}{cc} 2 & 3 \\ 3 & 5 \\ 6 & 6 \\ 1 & 8 \end{array}$
 $\bar{x} = 3 \quad \bar{y} = 5.5$

$$\text{Cov}(x,y) = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$(2-3)(3-5.5) + (3-3)(5-5.5) + (6-3)(6-5.5) + (1-3)(8-5.5) \div \frac{n-1}{4-1} = \frac{(-1)(2.5) + 0 + 3 \times 0.5 + (-2 \times 2.5)}{3} = -\frac{1}{3} \Rightarrow -0.33$$

What is the meaning of -0.33

It means 2 features (x, y) is very related.

$$= \text{Cov}(x,y) = +ve$$

$+ve$ -
 y, y is +vely related

very not talking about [Magnitude]

$$\text{Cov}(x,y) = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\text{Let } \text{Cov} = 0.5$$

Cov Sq. kg (we do this)
but Is it make any sense.

Name	Owner	Age	Sq. mtr	Wt/km	Time	Y
Jay	1800	15	10	1.2	20	1.225

If $\text{Cov}(f_1, f_3) \xrightarrow{\text{let}} \text{Sgft. years} = 0.6$

So we can make statements that

↳ The Strength of relationship

$$\text{b/w } f_1 \& f_3 \xrightarrow{0.6 > 0.5} f_1 \& f_2$$

→ unit same
Unit different

and magnitude is not helpful

X

Since the dimensions are different we can't compare the magnitude.

Then how
And Question 2 → What would be the value of $\text{Cov}(x, y)$ Answer: -∞ to ∞

① If $\text{cov}(f_1, f_2) = 0.1$ → then we say $\text{cov}(f_1, f_3)$ is twice $f_1 \& f_2$ → we say X

$$\text{cov}(f_1, f_2) = 0.2$$

↓

↳ why
because value is -∞ to ∞

→ Two Disadvantages of Covariance:

- ① Dimension/unit is there
- ② -∞ to +∞

→ Then what is use of Cov.

↳ takes about direction of relationship

↳ only

→ Then what will be solution → in Magnitude

↓ come & related

book

Shahrukh

Correlation

$$= \frac{\text{Cov}(x, y)}{\sigma_x \times \sigma_y} \rightarrow \text{eg. } \frac{\text{kg. sgft}}{\text{kg. sgft}} \rightarrow \text{So we say dimensionless quantity}$$

Solve 1 Problem

When divide x, y so the Cov range in -1 to 1

Solve Problem ② also

df. corr()

H.W Example \Rightarrow

X	Y
5	3
8	4
9	5
6	7
3	8

we need to know

⇒ How Correlations Works?

$$\text{eg.: } f_{x,y} = 0.4 \quad f_{A,B} = 0.8$$

↳ we say $f_{A,B}$ is highly correlated than $f_{x,y}$

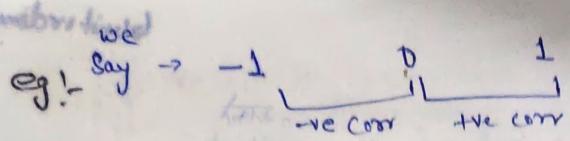
↳ $f_{x,y} = 0.4$ \rightarrow if 1 unit change in x then 0.4 time change in y features.

0.4
① +ve correlate

Correct Interpretation is

★ How much sure it is if there is 1 unit change in x the other unit will also change.

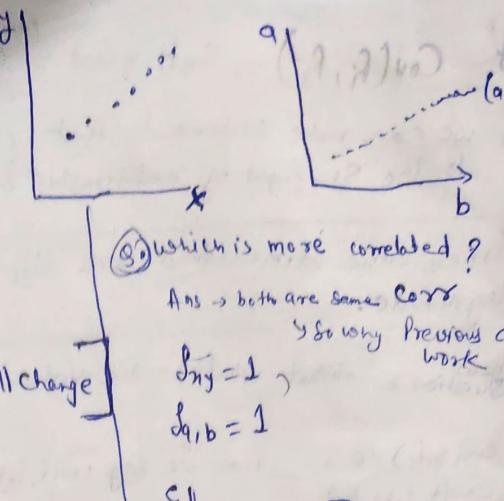
↳ note: 40% change may be possible in increase & decrease

eg! say \rightarrow 

→ Correlation is not about slope \hat{y}

So $\rightarrow 1 \Rightarrow$ perfect correlation

If 1 unit ↑ then 100% chance y will change

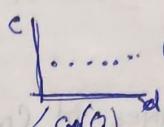


Ans \rightarrow both are same corr

↳ So why previous definition work

$$\hat{y}_{xy} = 1$$

$$d_{a,b} = 1$$



③ what is correlation

what tell these relationship (\rightarrow Spearman rank correlation)



④ Pearson Correlation Coefficient?

It always measures linear relationship

$$\text{Cov}(xy) = \frac{\text{Cov}(x,y)}{\sigma_x * \sigma_y}$$

Spearman rank Correlation \rightarrow Non-linear relationship.

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sqrt{R(x)} * \sqrt{R(y)}}$$

$R(x) = \text{Rank}(x)$

Spearman corr = 1
Pearson corr = 0.88

what Rank of x, y row by arrange or sort Value

x	y	$R(x)$	$R(y)$
5	6	3	1/2
7	4	2	3
8	3	1	5
1	1	5	5
2	2	4	4

Sort $\rightarrow \text{Rank}(x) = 1, 2, 5, 8, 7, 3$

~~6 5 4 3 2 1~~

$\text{Rank}(y) \Rightarrow 1, 2, 3, 4, 6$

~~2 3 4 5 1 6~~

Use Case:

31K features in house price prediction. \rightarrow How will you use correlation?

Any calc corr to see relation if

Close to 1 or -1 then use in ML

If $\text{Corr}(x) = 0$ drop that after consult BI.

in dataset now use
Spearman Correlating
always use
Pearson Correlate

→ Use heatmap

(good corr $(0.3 \leq 1)$ or $-0.3 \leq 1)$)

Probability

slightly

Random Variable

- ① Experiment → A repeatable procedure with a set of Possible result & eg tossing a coin
- ② Sample Space → All possible outcome of an experiments.
- Sample space → Head & tail.
- ③ Event → One or more outcome of Experiment.

* Basic Counting principle:-

→ m ways to do one thing and n way to do another
 Then there are total $m \times n$ ways to do both the things.

① Independent event
 ② Exclusive

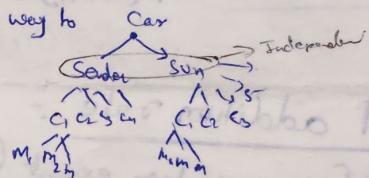
Eg: { 3 Shirts }
 { 4 jeans } → $3 \times 4 \rightarrow 12$ ways combi

② { 3 ice cream }
 { 4 cones } → $4 \times 3 \rightarrow 4 \times 3 = 12$

More than two choices

✓ Option 1 → Sedan or SUV
 Color opt 2 → C₁, C₂, C₃, C₄, C₅
 Opt 3. - Model - (M₁, M₂, M₃)

How many ways to buy car



$$2 \times 5 \times 3 = 30 \text{ ways}$$

• These basic counting principle work in independent events. P Adv

Probability principles:

• What is Probability?

Prob → Share of Success / Total no of Possible outcomes

Can we say //

It is Relative frequency

• Histogram
 • Scatter plot
 • P.D.F
 • Post
 • Dist. P.D.

Example: Prob Head → $P(H) = \frac{1}{2} \rightarrow$ Share of success
 We say say RF is what is share of Head → Total outcome

→ together no 2 will be 1

$$(8)^9 + (1)^9 = (8+1)^9 = 0$$

Principle

- ① The exact outcome cannot be predicted (exception). ↳ Shakuimama
- ② All possible outcomes are known.
- ③ Equally likely outcomes. ↳ means eg: if there is fair dr H & T chance equally out
- ④ Repeatable under uniform condition. ↳ means eg: if you fix in aged dice, friction, speed, rotation, then come only 1 out. eg G H but Ideally strictly not possible.

Probability Rules

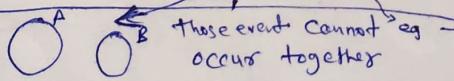
- ① For any event A, $0 \leq P(A) \leq 1$ ↳ why no > 1 because it is taken or you divide by total number
eg:
- ② Sum of all probability = 1. ↳ eg: toss a coin Success $\rightarrow (H, T)$ Total no on $(H, T) \rightarrow 2$
eg: $P(H) + P(T) = 1$
 $P(H) = 1 - P(T)$ {Rule of Subtraction}
- ③ Complement rule $\rightarrow P(\text{not } A) = 1 - P(A)$

④ General addition rule:

↳ for any two events

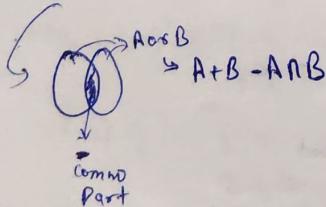
$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

Special case of disjoint events $\rightarrow P(A \text{ or } B) = P(A) + P(B)$



$P(A) + P(B)$ eg: Move you head in one go left & right
↳ travel by train & airplane same time

$$P(A \text{ or } B)$$



⑤ General multiplication rule:

$$\begin{aligned} P(A \text{ and } B) &= P(A) * P(B/A) \\ &\downarrow \\ &= P(B) * P(A/B) \end{aligned}$$

Special A and B are independent ↳ independent event \rightarrow does not affect each other may be occurs together

$$P(A \text{ and } B) \leftarrow P(A) * P(B)$$

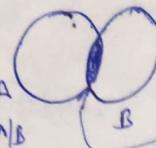
$P(A|B) \Rightarrow 0$ (Prob A when B has occurred.)

eg: I having breakfast & Ms Modi in Bill in LK 10pm
eg: I having breakfst & eating youtube, 10pm

Conditional Probability:

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

P(A) when event B is happen

 Balrady occur
जब अपने दो घटनाएँ होती हैं तो क्या
A घटना होती है तो वो क्या
Intersection part होती है आकर्षणीय
जो यह योग के B के

$$\frac{P(A \cap B)}{P(B)} \leftarrow \text{शुद्धिकरण प्रक्रिया के बारे में Balrady होती है}$$

$\frac{\text{शुद्धिकरण}}{\text{प्रक्रिया}}$

$$\leftarrow \text{Total probability}$$

Now go to Rule ④ for better understanding

Q. 16 people study french, 21 study spanish. There are 30 people altogether.

Work on probabilities

$$P(F) \rightarrow \frac{16}{30}$$

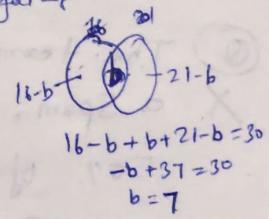
$$P(S) \rightarrow \frac{21}{30}$$

$$P(F \text{ only}) \rightarrow \frac{9}{30}$$

$$P(S \text{ only}) \rightarrow \frac{12}{30}$$

$$P(F \text{ or } S) \rightarrow \frac{27}{30}$$

$$P(F \text{ and } S) = \frac{(F \cap S)}{30} = \frac{7}{30}$$



Special case of condition

Bayes Theorem

$$\text{we know Conditional Probability} \Rightarrow P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{--- (1)}$$

$$\text{or} \quad P(B|A) = \frac{P(B \cap A)}{P(A)}$$



$A \cap B \Rightarrow$ are same yes
 $B \cap A$

$$\text{we say} \quad P(A \cap B) = P(A|B) * P(B) \quad \text{--- (2)}$$

$$P(B \cap A) = P(B|A) * P(A) \quad \text{--- (3)}$$

equality (1) & (2)

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

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

Bayes theorem

→ Naive Bayes
ML
→ Most classification

Revision

- ① A dice two prob 6? → 1/6
- ② A bag 5 red & 3 blue → Random. Prob of Yellow $\Rightarrow \frac{5}{16}$
Yellow ①

8 52 card \Rightarrow Prob hearts/Red? $\frac{13}{52} \Rightarrow \frac{1}{4}$

Likelihood \rightarrow (Prob of) reverse happening

Prior Prob

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

Posterior probability

Evidence

What is updated prob of A if some evidence of event B is there
thinkon $P(A|B)$ \rightarrow Its meaning that evidence of B given rd prior prob of A was $P(A)$.

wrong value

X not good even

③ The term free occurs in 20% of the emails marked as

X a Spam. 0.1% of non-spam emails include the term free.

50% of the emails are spam.

Calculate prob that an email is spam if the word free appears in it. $\Rightarrow P(\text{spam}/\text{free})$

$$\Rightarrow P(\text{spam}) = 0.5$$

$$P(\text{free}/\text{spam}) = 0.2$$

$$P(\text{free}/\text{nonspam}) = 0.01 \Rightarrow 0.01\%$$

$$P(\text{spam}/\text{free}) = \frac{P(\text{free}/\text{spam}) \cdot P(\text{spam})}{P(\text{free})}$$

$$= \frac{0.2 \times 0.5}{0.01} \Rightarrow \frac{0.10}{0.01} = 10\%$$

④ 10% of patients in a clinic have liver disease. 5% of the clinic patient are alcoholics. Among these patients diagnosed with liver disease 7% are alcoholics. You are interested in knowing the probability having liver disease given that he is an alcoholic.

$$P(A) = \text{Prob of having liver disease} = 0.10$$

$$P(B) = \text{Prob of alcoholic} = 0.05$$

$$P(B|A) = \text{having liver disease who has alcoholism} = 0.07$$

$$P(A|B) \rightarrow \text{Prob drink alcohol given alcohol patient having liver disease stat}$$

$$\frac{P(B|A) \cdot P(A)}{P(B)} \Rightarrow \frac{0.07 \times 0.10}{0.05}$$

$$\Rightarrow 0.14 = 14\%$$

Q2 In a factory 30% of the products are defective. If a defective product is randomly selected, what is the probability that it was produced by Machine A, given that Machine A produces 60% of defective products?

$$P(\text{defective}) = 30\% \\ P(D) = \frac{P(D|A) \cdot P(A) + P(D|B) \cdot P(B)}{P(D)} \rightarrow \text{Law of total Prob}$$

$$P(D) = P(D|A) \cdot P(A) + P(D|B) \cdot P(B)$$

Use Case!

Eg. ① $P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$ eg: House Price f_1, f_2, f_3 | Price
so we write like

$$P(f_1, f_2, f_3 | y) = \frac{P(y | f_1, f_2, f_3) \cdot P(f_1, f_2, f_3)}{P(y)}$$

$$= \frac{P(y | f_1) \cdot P(y | f_2) \cdot P(y | f_3) \cdot P(f_1) \cdot P(f_2) \cdot P(f_3)}{P(y)}$$

Eg. ② What is price if room area loyalty

Naive Bayes Model
In ML called N.

- AI
- Medical testing
- financial risk
- Marketers (feedback)
- Medical diagnosis

		Play/not play	outlook	temp	play
			Sunny	Hot	No
			Rainy	Cold	Yes
			overcast	Mild	
				Hot	
				Cold	
				Mild	

$$\begin{aligned} P(Y) &= 9/14 \\ P(N) &= 5/14 \end{aligned}$$

③ Today (Sunny/Hot) → Play or not

↙ so here use bayes theorem

$$P(\text{Yes} | \text{Today}) = \frac{P(\text{Sunny}/\text{Yes}) \cdot P(\text{Hot}/\text{Yes}) \cdot P(\text{Yes})}{P(\text{Today})}$$

68

$$P(\text{No}/\text{Today}) = \frac{P(\text{Sunny}/\text{No}) \cdot P(\text{Hot}/\text{No}) \cdot P(\text{No})}{P(\text{Today})} \rightarrow \text{constant for both same}$$

P(Y)	P(N)	P(Yes)	P(No)
2/9	3/5	2/9	2/5
4/9	1/5	4/9	1/5
3/9	2/5	3/9	2/5
100%	100%	100%	100%

P(Yes)	P(No)
2/9	2/5
4/9	1/5
3/9	2/5
100%	100%

$$P(\text{Yes}) = 0.031 \rightarrow \text{Normalize these} \\ P(\text{No}) = 0.085 \rightarrow \frac{P(\text{Yes})}{P(\text{Yes}) + P(\text{No})}, \frac{P(\text{No})}{P(\text{Yes}) + P(\text{No})}$$

$$\frac{0.03}{0.03 + 0.085} \rightarrow \frac{0.035}{0.03 + 0.085}$$

$P(X=0.5) \rightarrow$ we
 law of Large Num.
 (1000)
 converges to
 average
 50% → Head
 50% → tail

Random Variable

Start

① Tossing a Coin..

Experiment is Random $\rightarrow H, T$

the outcome will be Random

How Represent in Maths

$$X = \begin{cases} 0 & H \\ 1 & T \end{cases} \rightarrow \text{quantified each of the event}$$

X called
Random Variables

$$② X = \{1, 2, 3, 4, 6\}$$

↳ discrete \rightarrow this is discrete random Variable
 which
 This is discrete
 or Cont.

in. Alzebra
 $x+5=0$
 $x=-5$
 $X \rightarrow$ value is fixed
 ↓

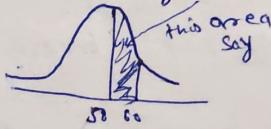
Experiment
 by
 random
 Outcome
 Random
 ↓
 discrete continuous

→ Continuous Random Variables

age
height

distribution of Age \rightarrow continuous \rightarrow then How to calculate the probability here

↳ ① if what is prob of age 50 & 60 Age



- A Random Var Can take any Values
- random Variable Value is unknown. $\{1, 2, 3, 4, 5, 6, \dots\}$

Eg: Toss $\{H, T\}$

$$\left. \begin{array}{l} P(H) = 1/2 \\ P(T) = 1/2 \end{array} \right\} \rightarrow \text{give me a generic function}$$

↳ give me a gf to calculate

Prob of a event while tossing a coin

$$P(x) = \frac{1}{n}$$

$n = \text{total no. of outcome}$

$$\begin{aligned} P(x=0) &= 1/2 \\ P(x=1) &= 1/2 \end{aligned}$$

(Eg. 2) dice $\{1, 2, 3, 4, 5, 6\}$ we say $\frac{1}{n}$ function that can be use

$$P(2) \Rightarrow P(1) = P(n) = \frac{1}{n} \rightarrow \text{to get probability}$$

$$P(3) \Rightarrow \frac{1}{6}$$

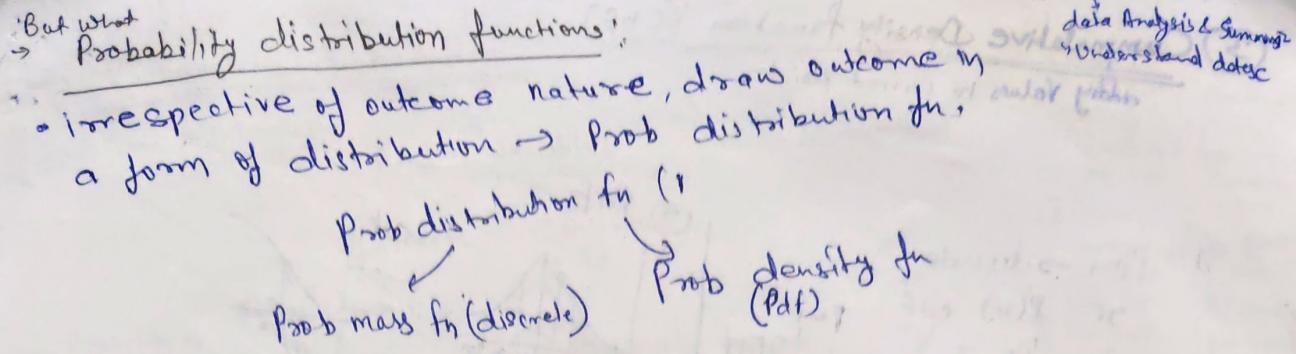
* Conclusion \rightarrow Outcome of Experiment

Discrete

Continuous

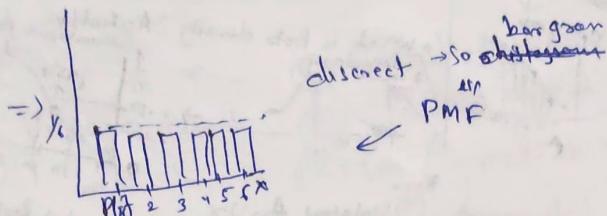
Prob mass function

Prob density fun

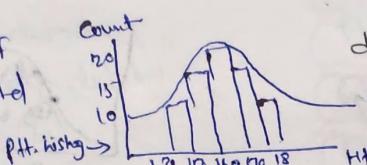


e.g.) throw dice

$$\begin{aligned} P(1) &= \frac{1}{6} \\ P(2) &= \frac{1}{6} \\ P(3) &= \frac{1}{6} \\ P(4) &= \frac{1}{6} \\ P(5) &= \frac{1}{6} \\ P(6) &= \frac{1}{6} \end{aligned}$$

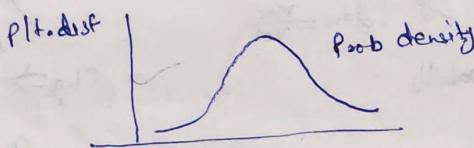


② Continuous - PDF
Height of Std.



PDF is a statistical term that describes the probability distribution of the continuous random var

- Distribution Use
↓
data Analysis & Summary
Understand data
↓
① Anomaly detection
↓
Food detection, spending
toys fall into behaviors
with fixed
② KDE (Kernel Density)
↓
Smith histogram
③ Bayesian Inference
↓
Update the DB of new evid



Deep dive

e.g. Rolling of dice \rightarrow PMF

$$P(x=1) = \frac{1}{6}$$

$$P(x=2) = \frac{1}{6}$$

$$\begin{aligned} P(x \leq 3) &= P(x=1) + P(x=2) + P(x=3) \\ &\Rightarrow \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} \end{aligned}$$

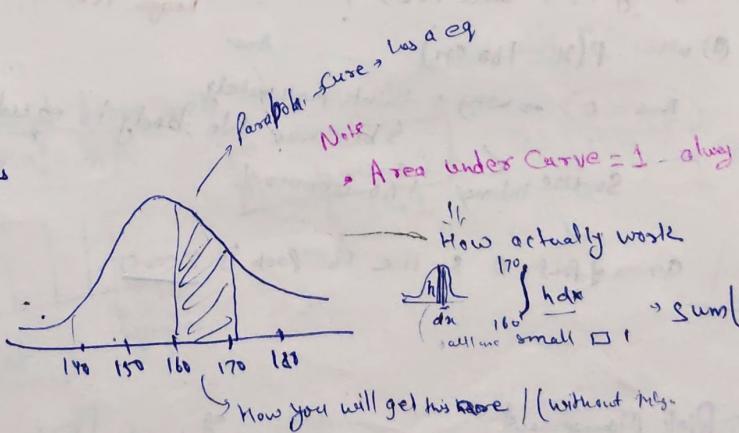
$$P(x \leq 4) = \frac{4}{6}$$

$$P(x \leq 6) = 1$$

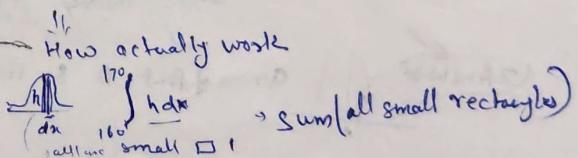
e.g. PDF \rightarrow continuous
Height of Std

③ What is Prob in Height
 $P(160 \leq X \leq 170)$?

or
18 ass prob(≤ 150)



Integration
use by Calculus
↓
Decide prob always 1



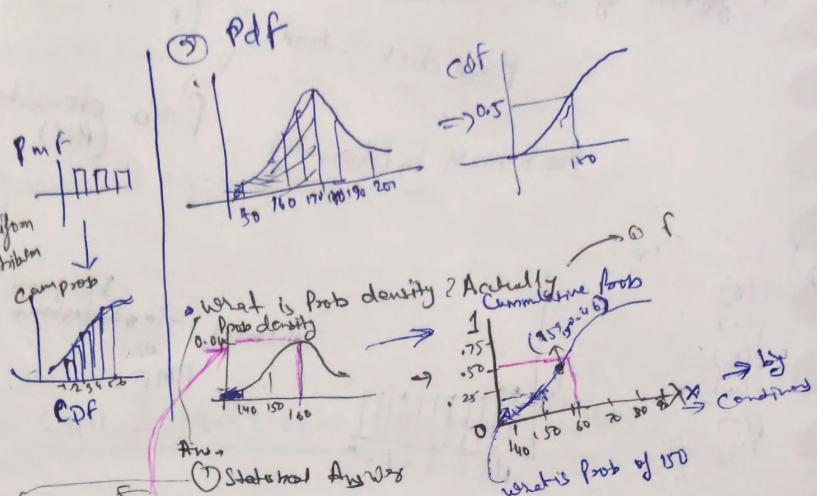
\rightarrow sum(all small rectangles)

③ Cumulative Density function (CDF)

adding values to current point

e.g.: Pmf \rightarrow discrete

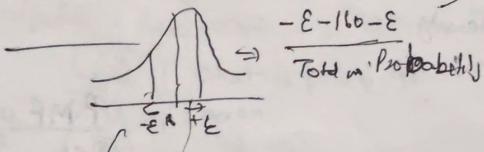
x	$P(x)$	CDF
1	$\frac{1}{6}$	$\frac{1}{6}$
2	$\frac{1}{6}$	$\frac{2}{6}$
3	$\frac{1}{6}$	$\frac{3}{6}$
4	$\frac{1}{6}$	$\frac{4}{6}$
5	$\frac{1}{6}$	$\frac{5}{6}$
6	$\frac{1}{6}$	$\frac{6}{6}$



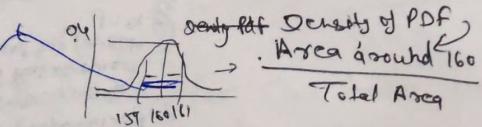
So we say that

\rightarrow Prob Density of Pdf
 \Leftrightarrow Slope of Cdf

PDF = Probabilty



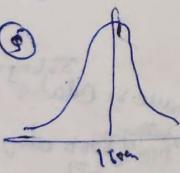
why



Note

Prob density can be greater than 1

Why took value of 159 and 161 ?



Q what $P(x=160 \text{ cm})$

Ans = 0

\rightarrow why 0 think mathematically
 \hookrightarrow because No body is equals to 160 cm

So use interval $160, 0.000001$

Q what $P(x>160)$

some one ask what is prob a person will have height 160cm

Ans

area of part = 0 so the 160 prob is zero.

2 Answer

\rightarrow USE Case \rightarrow ① Risk Management

Q calculate the prob of different risk level: (help determine the prob that an asset's value will fall below a certain threshold, which is critical for risk assessors)

② Quantile Calculation

\hookrightarrow Settly threshold for decision-making, ex: high-risk category

③ Hypothesis testing

\hookrightarrow calculate the p-value in statistical tests

e.g. A/B test

Probability distribution

we study it

→ Random Variables

Discrete

Continuous

If outcome is ω_j

then Probability = Use fn

probability mass
funet

distibn

- ① Discrete Uniform Discrete
- ② Bernoulli
- ③ Binomial
- ④ Poisson distribut
- ⑤ Geometric distribut

Probability function

J

- ① Normal distribution
- ② Standard normal distribution
- ③ Log-Normal distribution
- ④ Chi-Square dist
- ⑤ F dist
- ⑥ Exponential dist
- ⑦ Continuous uniform dist
- ⑧ t-distribution.

→ Prob function
PMF, PDF

→ What we see in
indistibn Distribut

- ① Notation
- ② Shape → & what
- ③ factor shaped by
- ④ pmf eq. opal
- ⑤ cdf eq. opal

⇒ In Uniform → Discrete

→ Continuous

① Discrete Uniform distribution ($U(a, b)$)

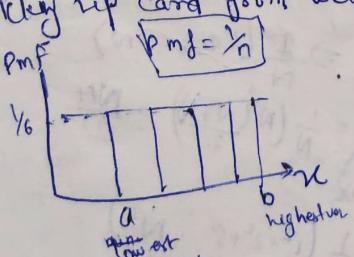
- discrete
- A uniform distribution refers to a type of probability distribution in which outcomes are equally likely.

→

Ex: Rolling a dice
 $\{1, 2, 3, 4, 5, 6\}$. The outcomes are independent

② tossing a coin

③ Picking up card from well



→ mathematically write

$$n = b - a + 1$$

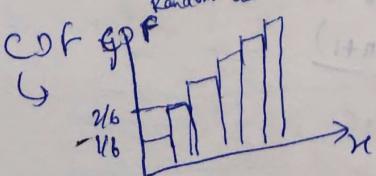
$$= 6 - 1 + 1$$

$$= 6$$

e.g.: $\{1, 2, 3, 4, 5, 6\}$

what is CDF?

Probability till the end



→ e.g.: As $1, 2, 3, 4, 5, 6$

for each of the distribution is mean
mean of discrete uniform distribution = $\frac{a+b}{2}$

Variance of discrete Uniform distribution = $\frac{n^3 - 1}{12}$

Here why we calculate the mean & Variance

we know mean & Variance give the information

Early learn → mean = $\frac{\text{sum of nos}}{\text{total no}}$ b/c

so related to mean there is terms/or concept

Expected Value → In the long run avg value of represented
of experiments

	x_i	$P(x_i)$	$x_i \cdot P(x_i)$
1	x_1	$P(x_1)$	$x_1 \cdot P(x_1) = 0.167 \times 1$
2	x_2	$P(x_2)$	$2 \times 0.167 = 0.167 \times 2$
3	x_3	$P(x_3)$	$3 \times 0.167 = 0.167 \times 3$
4	x_4	$P(x_4)$	$4 \times 0.167 = 0.167 \times 4$
5	x_5	$P(x_5)$	$5 \times 0.167 = 0.167 \times 5$
6	x_6	$P(x_6)$	$6 \times 0.167 = 0.167 \times 6$

mean of Uniform distib.
 $M = \frac{1+2+3+4+5+6}{6} = \frac{21}{6} = 3.5$

$\sum_{i=1}^n x_i \cdot P(x_i)$ around 3.5

what you observes here

So we can say that mean & the Expected Value is same

$E(x)_{\text{full}} = \sum_{i=1}^n x_i \cdot P(x_i)$

↳ Variance of Uniform distribution

$\text{Var}(x) = E(x^2) - \frac{M^2}{12}$

eg. $\text{Var}(x) = \frac{6^2 - 1}{12} = \frac{35}{12} = 3.083$

prove ↳

sum of n natural numbers = $\frac{n(n+1)}{2}$

sum of n^2 numbers = $\frac{n(n+1)(2n+1)}{6}$

we have $\text{Var}(x) = E(x^2) - (E(x))^2$
Variance of any Number

$E(x) = \sum_{i=1}^n x_i \cdot P(x_i)$

$= \sum_{i=1}^n x_i \cdot \frac{1}{n} \Rightarrow \frac{1}{n} (1, 2, 3, \dots, n)$

$= \frac{1}{n} \left(n \cdot \frac{(n+1)}{2} \right) = \frac{n(n+1)}{2}$

$E(x^2) = \sum_{i=1}^n x_i^2 \cdot P(x_i)$

$= n^2 \cdot \frac{1}{n} \Rightarrow \frac{1}{n} (1^2 + 2^2 + 3^2 + \dots + n^2)$

$= \frac{n(n+1)(2n+1)}{6}$

$\text{Var}(x) = E(x^2) - (E(x))^2$
 $\Rightarrow \frac{(n+1)(2n+1)}{6} - \left(\frac{n(n+1)}{2} \right)^2$

$\text{Var}(x) = \frac{n^2 - 1}{12}$

$\sigma = \sqrt{\frac{n^2 - 1}{12}}$
SD

Continuous Uniform distribution

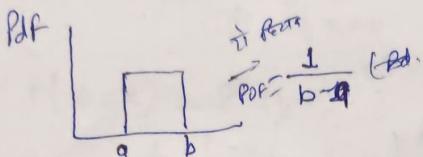
- That has infinite no. of values in a specific range
- R.V is continuous
- It is rectangular distribution

- Ex:
- ① A perfect random no. generator
 - ② prob of guessing exact time at any moment
 - ③ Waiting time at a bus stop
 - ④ Temporal variations in a day

① Notation: $U(a, b)$

Parameters: $-\infty < a < b < \infty$

②

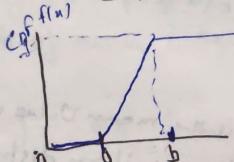


what is P that area $= b-a \times f(x)$ | we know area under give the prob which is 1

$$1 = b-a \times f(x)$$

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

How to make CDF



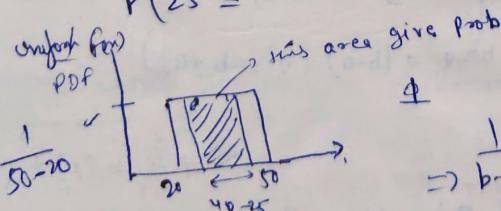
$$\text{CDF} = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & x \in [a, b] \\ 1 & x > b \end{cases}$$

→ So we see prob lie in CDF/PDF & a & b is start & end

→ The mean of Cont Uniform dist = $\frac{1}{2}(a+b)$
Variance: " " " " = $\frac{1}{12}(b-a)^2$

③ The no. of items sold at a shop daily is Uniformly distributed with max and min sold 50 & 20 respectively so what is prob of daily sales to fall b/w 25 to 40 ??

$$P(25 \leq x \leq 40) = ??$$



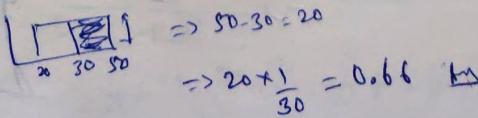
1

$$\Rightarrow \frac{1}{b-a} \quad \text{height start end } 50-20$$

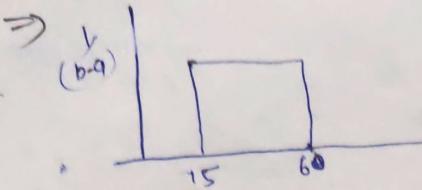
$$\frac{1}{50-20} = \frac{1}{30}$$

$$\text{prob}(f(x) \leq v) = \frac{1}{30} \times 15 = 0.5$$

$$P(x \geq 30) \rightarrow$$



Q. The amount of time for pizza delivery is Uniformly distributed b/w 15 & 60 mins. What is standard deviation of the amount of time it takes for pizza to be delivered?



$$\text{Var} = \frac{1}{12}(b-a)^2$$

$$= \frac{1}{12}(60-15)^2 = \frac{(45)^2}{12} = 168.75$$

$$SD(x) = \sqrt{\text{Var}} = \sqrt{168.75} \approx 13 \text{ min}$$

Optional

→ for a continuous random var prob density fn ($f(x)$)

$$\rightarrow \text{by defn } E(x) = \int_a^b x \cdot f(x) dx$$

\rightarrow Sum integral
Summation for contn variable

\rightarrow before for Discrete
 $\sum x_i \cdot p(x_i)$

$$\hookrightarrow \text{Uniform distribution } \int_a^b x \cdot f(x) dx \Rightarrow \int_a^b x \cdot \frac{1}{b-a} dx \Rightarrow \int_a^b x - \frac{1}{b-a} dx$$

Basic Integration formula

$$\int x^n dx \Rightarrow$$

$$\frac{x^{n+1}}{n+1} \Rightarrow \frac{x^2}{2}$$

$$\Rightarrow \frac{1}{b-a} \int_a^b x dx \Rightarrow \frac{1}{b-a} \left[\frac{x^2}{2} \right]_a^b \Rightarrow \frac{1}{2(b-a)} [x^2]_a^b$$

$$\Rightarrow \frac{1}{2(b-a)} (b^2 - a^2) \Rightarrow \frac{1}{2(b-a)} (b+a)$$

↑ upper value

$$\Rightarrow \boxed{\frac{b+a}{2}}$$

this is the mean value of
Continuous uniform distribution

⇒ Variance of any continuous distribution

→ It help to derive any variance of distribution

$$\text{Var}(x) = E(x^2) - (E(x))^2$$

$$E(x^2) \Rightarrow \int_a^b x^2 \cdot f(x) dx \Rightarrow \int_a^b x^2 \cdot \frac{1}{b-a} dx \Rightarrow \frac{1}{b-a} \int_a^b x^2 dx$$

$$\Rightarrow \frac{1}{b-a} \left[\frac{x^3}{3} \right]_a^b \Rightarrow \frac{1}{3(b-a)} (b^3 - a^3)$$

$$\Rightarrow \frac{1}{3(b-a)} (b^2 + ab + a^2)$$

$b^3 - a^3 = (b-a)(b^2 + ab + a^2)$

$$= b^2 + ab + a^2$$

$$\text{Put } = E(x^2) - \frac{1}{3}(E(x))^2 \Rightarrow \left(\frac{b^2 + ab + a^2}{3} \right) - \left(\frac{a+b}{2} \right)^2 \Rightarrow \frac{b^2 + ab + a^2}{3} - \frac{a^2 + b^2 - 2ab}{4}$$

$$\Rightarrow \frac{4b^2 + 4ab + 4a^2 - 3a^2 - 6ab - 3b^2}{12} \Rightarrow \frac{b^2 - ab + a^2}{12}$$

$$\Rightarrow \boxed{\frac{(b-a)^2}{12}}$$

If it is a
(Variance distn)
of the Uniform
distribution

Types Prob Distrn

② Bernoulli distribution

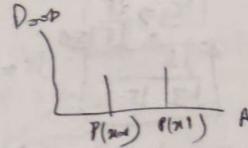
A discrete Prob dist. of a random var which takes only two possible outcomes typically labelled as success (1) and failure (0)

To model a problem statement with only two possible outcome

e.g.: Toss a coin (head or tail)

$$\begin{cases} P(X=0) = 0.5 = p \\ P(X=1) = 1 - p = 1 - 0.5 = q \end{cases}$$

write Mathematically



$$PMF \Rightarrow P(X=k) = \begin{cases} p & \text{if } k=1 \\ 1-p & \text{if } k=0 \end{cases}$$

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

$$\begin{aligned} \text{if } k=1 &\rightarrow p^1 (1-p)^{1-1} \\ &\Rightarrow p^1 (1-p)^0 \\ &\Rightarrow p^1 \cdot 1 \Rightarrow p \\ \text{if } k=0 &\rightarrow p^0 (1-p)^{1-0} \\ &\Rightarrow 1 \cdot (1-p) \Rightarrow (1-p) \end{aligned}$$

\Rightarrow Condition of BD

- ① No of trial should be finite
- ② Each trial should be independent
- ③ Only two possible outcome
- ④ Prob of each output should be same in every trial

Example:

- ① tossing a coin (binomial or Discrete dist.)

② prob of getting [3 and not 3] while throwing dice

In ML.

In classification

So we use Bern		dataset	target var	count	loan	salary	credit/default
1000	1000	1	0	2	5	3	1
1000	1000	2	1	0	0	0	0

for this classificat Prob
get a cost for which
is actually comp from
bernoulli \rightarrow PMF

Cg: defaulter/not defaulter

Pass/fail in exam

into logistic regression
using Bernoulli

\Rightarrow Mean and Variance of Bernoulli distribution

$$\# \text{Mean } E(X) = \sum_{i=1}^k x_i p(x_i) \quad \begin{aligned} &\text{let } p=0.6 \\ &\Rightarrow X=1 + 2 \cdot 0 \end{aligned} \quad \begin{aligned} &\text{let } 0.4 \\ &\Rightarrow X=0 + 1 \cdot 0.4 \end{aligned}$$

$$\# \text{Var } E(X^2) - E(X)^2 \rightarrow (p)^2 \quad \Rightarrow 1 \cdot 0.6 + 0 \cdot 0.4 = 0.6 \Rightarrow p$$

$$\epsilon (x)^2 = \sum x^2 (P(x)) \rightarrow (1-p)$$

$$\Rightarrow 1^2 \cdot p + 0^2 \cdot (1-p)$$

$$\Rightarrow p = \frac{E(X^2) - (E(X))^2}{(1-p)} \Rightarrow p - (1-p)^2 = p(1-p) \text{ proved}$$

③ Binomial distribution

$2 \rightarrow 2$ outcome

- Binomial distribution is ' n ' \rightarrow Bernoulli trials

So $P(X=k) \rightarrow$

$$n \choose k p^k (1-p)^{n-k}$$

↑
no of trials

for Bernoulli
 $p^k (1-p)^{1-k}$

How to Interpret

↓

e.g.: - 1st toss 2nd, ... 10th toss

What's Bay binomial

↳ What is prob of getting 'k' heads out of n trials

$n \rightarrow$ total no of trials

$k \rightarrow$ the no of events that you are get.

e.g. with 3 tosses what is prob of getting exactly 2 head ??

$${}^n C_k p^k (1-p)^{n-k}$$

$$\Rightarrow n=3, k=2 \quad \text{prob of head} = 1/2$$

$$\cdot P(X=2) = {}^n C_k p^k (1-p)^{n-k}$$

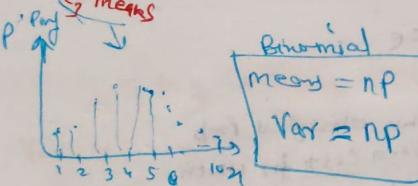
$$\Rightarrow {}^3 C_2 \left(\frac{1}{2}\right)^2 \left(1-\frac{1}{2}\right)^{3-2} = \frac{1}{2} \cdot \frac{1}{2} \cdot (0.5)^2 (0.5)$$

$$\Rightarrow \frac{3 \times 2}{2 \times 1} \cdot (0.5)^3 = 3 \times (0.5)^3 = 3 \times 0.125 = 0.375$$

④ What you toss a coin 10 times what is prob that you will get head exactly 3 times

$$\Rightarrow n=10, k=3, \quad \text{Prob} = 1/2$$

what is 10 means



⑤ Po.

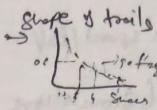
④ Poisson distribution

④ Geometric distribution

↳ It represents the probability of getting first success after having a consecutive no of failure (after n no of trials)

e.g.: I am interested in getting (H)
1st toss T 2nd toss T 3rd toss T 4th toss H

↳ you interest and now stop the trail → this is called as G.D.



$$\text{PMF} \Rightarrow P(X=x) = (1-p)^{x-1} \cdot p$$

e.g. G.D
jumping ball → Min. →

→ Not so much use case

$$\begin{aligned}\text{Mean} &= np \\ \text{Var} &= \frac{1-p}{p^2}\end{aligned}$$

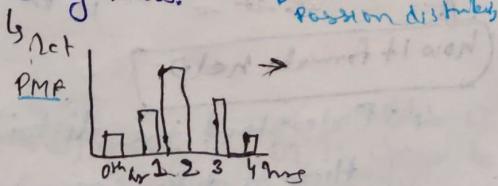
③ Poisson distribution

It is a discrete prob distribution that describes the no of events that occur within a fixed interval of time or space given a known average rate of occurrence.

→ important to BS, data analyst,

e.g.: No of events occurring in a fixed time interval.

① No of calls received by a customer every hour.
Expected no of calls every hr is $\lambda = 100$



↳

② No of people going to temple
... to hospital/airport every hour

$$\text{PMF} = P(X=x) = \frac{e^{-\lambda} \cdot \lambda^x}{x!}$$

$\lambda = \text{avg. rate of events every hr/interval}$
 $\lambda = 2.718$

③ If $\lambda = 10$

Prob of a person visiting at 5th hr

$$P(X=5) = \frac{e^{-\lambda} \cdot \lambda^5}{5!} =$$

④ The avg no of customers entering a store in an hour is 5
What is the prob exactly 3 customers will enter the store next hour

$$\rightarrow \lambda = 5 \quad P(X=3) \Rightarrow \frac{e^{-\lambda} \cdot \lambda^3}{3!} \Rightarrow \frac{e^{-5} \cdot 5^3}{3!} \Rightarrow \frac{(2.718)^5 \times 125}{6} \Rightarrow \frac{0.00674 \times 125}{6} \approx 0.214$$

As 14% chance to come 3 persons

No. of Probability density function (for continuous random variable)

① Normal distribution | Gaussian distribution | Bell-shape distribution

why most see in data

↳ Most of the naturally occurring/human generated follow a normal distributions.

↳

Why

Example ① Buincunx (math is fun)

Surely → ↗.

→ symmetry
world

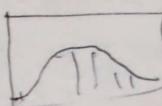
no matter what is speed / ball

② Galofton board (simulations)

What is point g show

Let again →

Experiment - Measure of height



→ The distribution helps in hypothesis testing
how will you do HT if you don't know what is the nature of distributions.

Why learn:

↳ ① Natural data follows -

↳ ② Statistically → by its characteristics

① Symmetric in nature

② Skewness is zero

③ Kurtosis = 3 (most best)

④ mean = median = mode

but not so much in each
bcz also find there of the

③ - 68-95-99.7% empirical rule.

(very empirical)

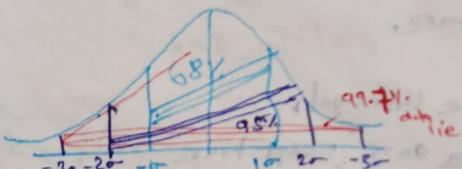
so it is true & that's why

with experiment / also

no any statistical proof also

bcz it is not proveable

like. lot of experiment



how it formula help?

It helps to understand the distribution of data

④ Outliers of data (→ how) → Box Plot

we say

extreme values (low, high) loss in position

young use

B, H, B

→ loss in position (normal data)

How say most data lie

in Q₁ to Q₃ ← come

we know

-50 to 10 (H, B)

or high -

or low -

beyond -3σ to 3σ

means ration or majority

or majorities

⑤ Prob about prob density

⑥ Calculate the prob

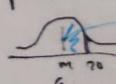
people will fall b/w 60kg & 70kg

→ how find

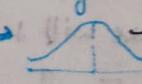
→ by calculate the area (AUC)

not manually

by use table



⑦ What is Prob if the wt is 60kg



ans is prob is 0, because the AUC of point is zero

Real Example: Ht weight, Mean of error
Score, Blood Pressure,

Can we mean = 0 } → for normal distribution
std = 1 } → strong

if we say it for
Standard Normal
Useless

→ Industry use case of normal distribution

↳ Beyond 3rd std → outliers

↳ Some of ML/Statistical model → require the data to be normally distributed → why → (1) Most of data fall natural

(2) Outlier identification (so deal with outliers)

(3) Prediction will be better

If data is skewed or not normal
No data will be dist

or get
Prediction

Something in Skewed data
Prediction will be wrong
(Under or Overfitted)

↳ why

→ If like this data → use transformation in data

Blog (feature Treatment)
inds.

1. Log transformation → for right-skewed data
2. Square → data \rightarrow log of data

- reciprocal
- Boxcox

and vs

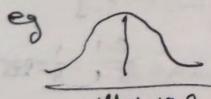
② Which transformation we when?

• There is rule → but also do hit & trial

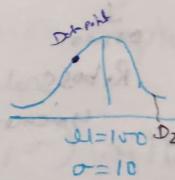
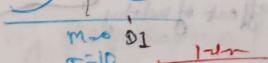
→ If you are modeling statistical ML models follow normal distribution and if not do transformation

mean = 0, Std = 1 → not

mean \neq 0, Std \neq 1 → center of



eg (1)



{ What is the $\sigma = 20$ } → interpolation

↳ On avg the data point is 20 unit away from mean

If we can't say

why → Both dist scale is not same. by size

Let say (2)

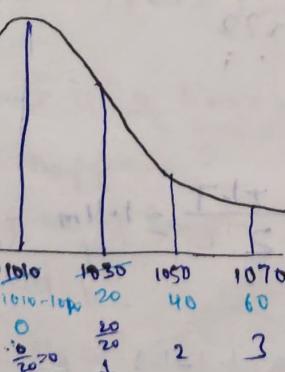
$$\mu = 1010$$

$$\sigma = 20$$

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

$$= \frac{1030 - 1010}{20}$$

$$= 1$$



Step 1

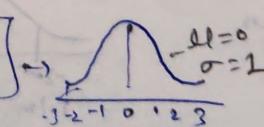
→ Can we subtract all with μ

If

part = $x - \mu$ has problem

Step 2 If we then divide by Std

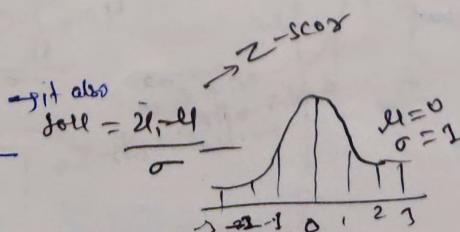
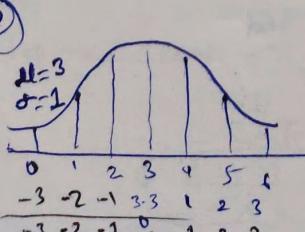
$$\frac{x - \mu}{\sigma}$$



Now we compare easily

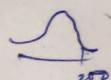
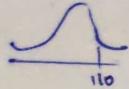
and it is called Standardization of data → $\mu = 0, \sigma = 1$

↳ Standard normal distn



Colored

⑧



1:53

Q Which Z-score is higher that deposit will be far from Mean

2 Industry use case of Std N.R.D

$$\text{Standardization} = \text{Z-score} = \frac{x_i - \mu}{\sigma}$$

Eg. house price Dataset

Room No.	Area of Price
1	2000
2	900
3	3000
4	1500

Worthy also

Both are in different Scale

→ unit is different also

→ for ML no problem in Prediction

→ but what is problem

do Scaling why.

↳ Better Implement

→ faster optimizat

(Gradient descent)

↳ Use by Scale factor

When we use which outlier has std deviation → remove

Outlier has std deviation → remove

remove effect of outliers

↳ others Robust Scales

Scaling (Optional)

↳ for small range

① Standardization ($\mu=0, \sigma=1$)

② Normalization (Range 0 to 1)

Min-Max Scaler

$$(x_{\text{norm}} = \frac{x_i - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}})$$

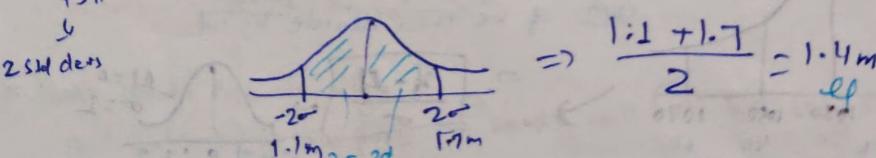
$$\begin{aligned} & \rightarrow \frac{1-1}{5-1} = 0 \\ & \frac{2-1}{5-1} = \frac{1}{4} = 0.25 \end{aligned}$$

2:25

Q 95% of Students at School are b/w 1.1m and 1.7m tall.

Can you calculate mean and std dev??

95% → 1.1m and 1.7m tall



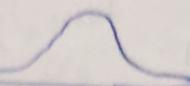
$$\Rightarrow \frac{1.1 + 1.7}{2} = 1.4 \text{ m}$$

$$4 \text{ std dev} \Rightarrow 1.7 - 1.1 \Rightarrow \frac{1.7 - 1.1}{4} = 0.15 \text{ m}$$



Central Limit theorem

(8th working of world)

We know the \rightarrow Normal Distribution:

Now power
ful
normal
empirical
family

 \rightarrow Symmetric \rightarrow Skewness ≈ 0 \rightarrow Kurtosis ≈ 3 \rightarrow mean = median = mode \rightarrow 67.95-99.7 rule

if distribution like



Haphazard/irregular Distribution \leftarrow We cannot define statistic
for irregulars!

\downarrow not useful for us. \rightarrow by statistical find away

CLT \rightarrow A way to convert any form of distribution to
Normal distribution.

\checkmark How works / or CLT says \rightarrow

$$\begin{aligned} \text{eg. } & S_1(50) \rightarrow \bar{x}_1 \\ & S_2(50) - \bar{x}_2 \rightarrow \\ & S_{50k}(50) - \bar{x}_{50k} \end{aligned}$$

Sample with replacement

\hookrightarrow a datapoint selected
in $S_1 \rightarrow$ may or may not
in S_2

$\begin{matrix} 1,2 \\ 3,4 \\ 5,6 \\ 7,8 \end{matrix}$

$s_1(1,2)$
 $s_2(3,4)$
 $s_3(2,4)$
 $s_4(3,4)$

\rightarrow The CLT states that if you have any pop with μ, σ and you take a sufficiently large random sample from population with replacement normally distributed.

- Sampling mean of population (μ, σ) will approximate to a normal distribution

Population $(\mu, \sigma) \rightarrow$ largely $\xrightarrow{\text{Sample}} \rightarrow$ Sample \rightarrow ND

$$\mu = N$$

$$\sigma = \sigma$$

$$\sqrt{n}$$

Question: (1) What is n here?
 \hookrightarrow Sample size $\text{e.g. } n=20$

\hookrightarrow in Sample Size 50

(2) why it happens?

(3) How many no. of sampling I take? \rightarrow Should be large

(4) What should be value of n ? ≥ 30 \rightarrow ND

(5) What is the usecase.

\rightarrow Any $n \geq 30 \Rightarrow$ sampling mean follow Z-distribution

Pop \Rightarrow Z-score $= \frac{x - \mu}{\sigma}$

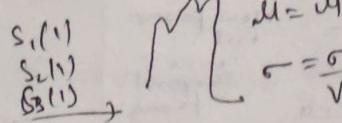
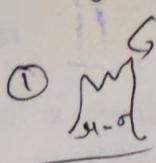
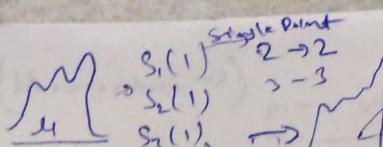
CLT \rightarrow Z-distribution

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

\rightarrow Z-distribution

<math

Why all is all after CLT?
Why σ is σ/\sqrt{n} after CLT?



$$\mu = 14 \quad \sigma = \frac{\sigma}{\sqrt{n}} \rightarrow \text{standard deviation}$$

$$-\frac{\sigma}{\sqrt{n}} \leq 0 \geq \frac{\sigma}{\sqrt{n}}$$

② $n =$

$S(\infty)$
whatever no of elements are there is
population you take all element in sample

$$\text{std} = 0 \quad \sigma = \frac{\sigma}{\sqrt{n}} = \frac{\sigma}{\infty} = 0$$

③ You have a population with $\mu = 100$ and std dev $\sigma = 20$

If you have sample size of 50 from this population,

What is the prob that sample mean will be less

than 105?

$$\Rightarrow \mu = 100, \sigma = 20, n = 50 \quad \bar{x} = 105$$

$$Z_{\text{score}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{105 - 100}{20/\sqrt{50}} =$$



$$\Rightarrow \frac{5\sqrt{2}}{4} = 1.7625 \quad \text{base } Z_{\text{score}} = 0.89$$

→ It's on line net
or sketch/CLT/ctb

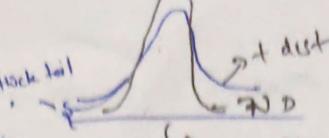
t-distribution: $\sim t(N)$

Student t-distrib

① Pop not given *

② What if Sample size < 30 ?

Student t-distribution



↳ many things distribute Look ND

ND

t-distribution *

• log-normal dist

• Exponential dist

• Power dist

• Chi-Square *

• F distribut *

• Pareto

$$\rightarrow \text{dof}(\text{D}) = n - 1$$

($D = +\infty$)

↳ becomes SND

$$\begin{aligned} \text{Z statistic} &= \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \\ \text{z-score} &= \end{aligned}$$

$$\begin{aligned} \text{T-Statistics Score} &= \frac{\bar{x} - \mu}{S/\sqrt{n}} \\ &= \end{aligned}$$

↳ Case \rightarrow T-test

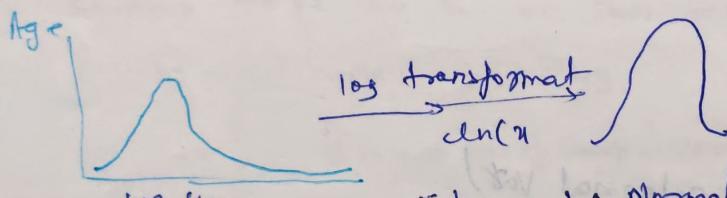
Intutions

• due to Sample Size

• Thick tail means the data is present at extreme

• log Normal distribut: A continuous distⁿ whose logarithm is Normally distributed.

$$x \rightarrow \log(x)$$



$\rightarrow \log \Rightarrow$ log Normal distribution

• used

example wealth distribution of world (rich others)

Use Case: \rightarrow \log

↳ where outliers can not be dropped, yet we use \log . N.D to make it a N.D

↳ isn't real & not be drop those as outlier

e.g.: time spent in reading Articles (1st part.....)

↳ people write comments on your post

Powers law dist (Pareto Principle)

Pareto distribution

rule

80/20 rule

20 - 11 times \rightarrow 80
called power law

e.g.: Party to Party = 100
20% choose
80% remain

20% employee

80% sell by

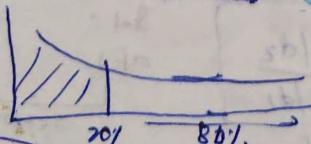
20% customer

80% of Rich 20 people

\rightarrow 80 thing done by top 20% people

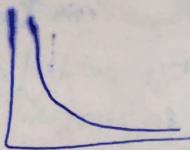
• Iple \rightarrow 2-3 player

• belts



$$\text{Pdf} = \alpha n_m^{\alpha} \quad \text{alpha parameter}$$

④ Exponential distribution:



$$\text{Pdf} = h e^{-hx} \quad h = .$$

- g. radioactive decay
- life span of battery

power law
Exp Pareto

⑤ CHI-SQUARE DISTRIBUTION: The CS dist is a prob that describes the distribution of sum of squares of K random variables.

Degrees of freedom

Population

Sample (2)

$S_1(1,2)$

$S_2(2,3)$

$S_3(4,5)$

$K = 2-1 = 1$

$\rightarrow S = \frac{(2-1)^2}{(3-1)^2} + \frac{(6-8)^2}{(5-7)^2}$

$\rightarrow f_{n,k}(x)$

$k=1$

⑥ Why always the CHI

→ Use case

① in ML (to compare 2 categorical vars)

② In hypothesis testing

⑦ f-distribution (fisher Snedecor distribution)

The f dist is a probability dist that is useful in context of comparing variables of two or more samples

→ It is right skewed takes only +ve values

→ F dist d_1 & d_2 is dist given by

$$F = \frac{S_1/d_1}{S_2/d_2}$$

Use case: compare two sample variance using f-test

S_1	S_2
Std -	Std
d_1	d_2

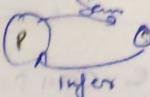
ratio
↳ ratio of
called Variance
F distribution

230 [23 March]

Inferential Statistics

↳ 6 ways estimate something about population using sample.

e.g.: Avg Salary of IT People in India is 15Lakhs



we are inferring something about pop' parameters using sample statistic.

⑥ Avg Salary of IT Employee is 15 Lakh.

↓ Statement

↳ Claim (Hypothesis)

define

• Hypothesis is a claim/predictive/generic statement, capable of being tested by scientific method.

eg: ① Avg age of std in a class of IIT is 23 yrs

② Consumption of Ice-cream sales increase in summers

③ Student who receive consulting will show greater increase in performance.

✗ ④ There is tone of water on Pluto.

✗ ⑤ Swargukt(hanson) we will go after we die.

✗ ⑥ In Patel dot stresh Naag lives.

⑧ How much this inference is true??



- You will pose some doubt
- You will avoid saying a specific no. > 15Lakhs

↳ Generation

↳ Sampling

↓

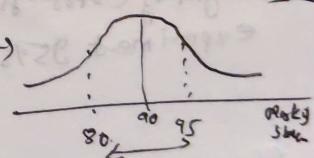
Point Est.

① If 89.95% < 100% No Conf.

② Interv → 80-95% more Conf.

↳ Inference

actually HT all about studying how confident you are with this particular inference



180-50
65-60
174.5
203.0

$$\text{Final} = 52.5 \pm 5$$

52.5 ± 5 ↑ Conf
52.5 ± 6 ↑ Conf
52.5 ± 7 ↑ Conf

Hypothesis Testing :

⑥ Avg Salary of IT Employee is 15 Lakh.

↓ Statement

↳ Claim (Hypothesis)

define

• Hypothesis is a claim/predictive/generic statement, capable of being tested by scientific method.

eg: ① Avg age of std in a class of IIT is 23 yrs

② Consumption of Ice-cream sales increase in summers

③ Student who receive consulting will show greater increase in performance.

✗ ④ There is tone of water on Pluto.

✗ ⑤ Swargukt(hanson) we will go after we die.

✗ ⑥ In Patel dot stresh Naag lives.

Example: Sachin: Avg age of Deloitte employee is 45 yrs.

test Alok → Sample 20 people → Sample Stats $\bar{x} = 49$ yrs

Agrim → 20 pe → Sample Stats = 42 yrs

★ Ruchika → 200 people S.S = 49.
Consider

↳ Highest no. of sample = 2000,
 $\bar{x} = 49$ yrs

S.E = 1.5 yrs

• framing of Hypothesis

Null H (H_0) = H_0 is a statement that is hold true unless we have sufficient statistical evidence to reject it.

Alternate (H_A): It is opposite of Null Hypothesis

H_0 : Avg Age = 45

H_A : Avg Age Not = 45

⑥ Avg Age of people in deloitte is greater than 45 yrs

Age > 45, Age ≤ 45 yrs

H_0 : Age > 45 H_A : Age ≤ 45

H_0 : Age ≤ 45 yrs, H_A : Age > 45 yrs

Whatever statement has equal sign that will be your Null Hypothesis

↳ Why

greater than 45
 H_0 (45 not included)

So we need to include the central point in order disprove something

Overall Inference

→ General Practice we care
G → Social Trial!

Level of Significance & Confidence!

→ We can not be 100% sure (Confident) to Conclude something for the population using sample.

not confident
the chances
of errors

e.g. If a pharma company wants to know if a vaccine will work or not.

→ 1000 times
95% → People well
↓ So
Vaccine works 95% of the time.

Out of 100 times, if I conduct a test

95% → make some threshold

how many times same type of conclusion can be made.

then what is

5%

→ How much margin of error you are ready to accept.

⇒ 5% → level of significance

So if I say I am ready to accept 5% margin of error then it simply out of 100 experiment 95%. If the true

The percentage of risk involved in testing. alpha → defined by business team

then the confidence is

$1 - \alpha$ → Confidence

$$\begin{aligned} \alpha &= 1,5\% \\ (I.) &= 99\% \\ &95\% \\ &90\% \end{aligned}$$

Steps

① Forming of Hypothesis

② alpha (margin of error)

③ What type of test that you will do?

Sample Statistics / Test Statistic } → Z-test
(Z, T, F, Chi-Sq)

or pop → SND } → Z-test

② S.S < 30 or → pop not

↳ + dist. F → t-test

③ Variance of two samples / or more than two
↳ F-distribution

④ Test Critical

↳ Z critical

↳ + critical

↳ f ??

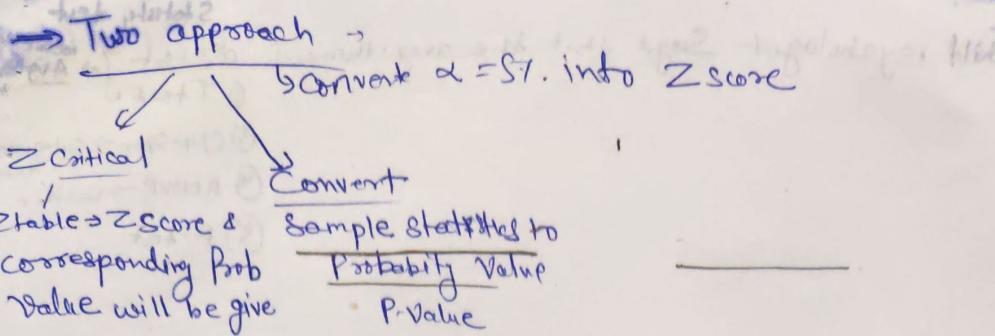
↳ Chi-Sq ??

Eg Sachin's H₀ slope = 45 yrs } 80 were rejected
H_A slope ≠ 45 yrs } SS = 49
we made for reject

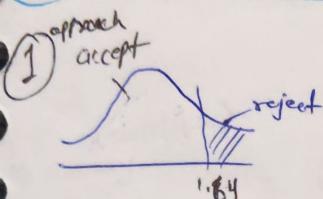
④ Two Categorical } → f test, ANOVA

↳ CHI-Square dist.

↳ Chi-Sq test



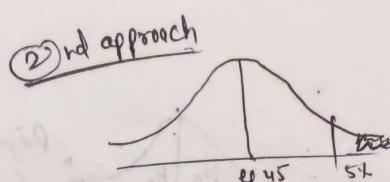
(5) Step 5 Conclusion



$$S.S > 1.91 \Rightarrow \text{res} \quad 5\%$$

$$\text{Statistics table} = 5\% = 1.64$$

$S.S > 1.64 \Rightarrow \text{reject region}$
you reject the H_0



↳ Prob value Correspondy
to Sample Score

$$P\text{-Value} \leq 0.05$$

↓
reject in H_0

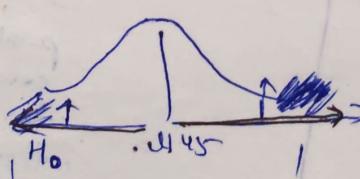
→ P-Value: The P-Value is the probability value, calculated from a test statistics

→ Use to decide whether to reject a H_0 or not.

*One tail test, two tail test:

$$H_0: \bar{x}_{age} = 45$$

$$H_A: \bar{x}_{age} \neq 45$$



$\neq \rightarrow$ Alternate Hypothesis of two tail test

↳ two tail test

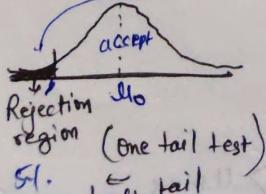
↳ alpha/2

↳ whole

Scenario-1 one tail

$$H_0: \bar{x}_x = \bar{x}_0$$

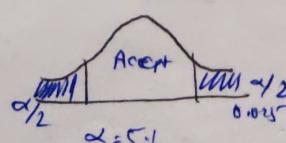
$$H_1: \bar{x}_x \leq \bar{x}_0$$



Scenario-2 Two Tail

$$H_0: \bar{x}_x = \bar{x}_0$$

$$H_A: \bar{x}_x \neq \bar{x}_0$$



Scenario-3 1 tail

↳ (Right tail)

$$H_0: \bar{x}_x = \bar{x}_0$$

$$H_1: \bar{x}_1 > \bar{x}_0$$



24/09
02:04

Q) Suppose a child psychologist says that the avg. time

Statistical test
Z-test \rightarrow Avg Var
T-test \rightarrow mean

- (3) CHI-SQR \rightarrow Categorical data
- (4) ANOVA \rightarrow Variance
- (5) F-test

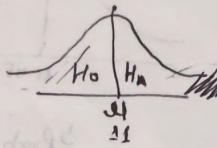
\rightarrow Step 1: $H_0: \bar{X} \leq 11$

$H_A: \bar{X} > 11$

Step 2 - Log Sig = $\alpha = 5\% = 0.05$

Step 3 \rightarrow S.S > 30 & pop give

$$\text{Z-test} \rightarrow \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \Rightarrow \frac{11.5 - 11}{2.3/\sqrt{100}} = 2.17$$



Rejection region
(one-sided)
dist. as sig.
1 tail test)

Step 4

(1) Z critical (2) P-value

- test Critical

$\alpha = 0.05$

Z table \rightarrow using α value

$\therefore \alpha = 0.05$

B6 $\leftarrow 0.0505$

$\therefore 0.05$

$\therefore 0.04$

3:13

② 2) Sechin

3:15 $x \approx N(4, 2)$

Q How many ~~employed~~ ~~not employed~~ or just in some I work such that no fit?

③ 3:22 Q what percentage

3) staff employed. Not with wage ab 1-7 earn 8%
what about

④ 3:28

⑤ 3:31

comes to specific financial condition

-> starting

Larval \rightarrow adult
young \rightarrow young

comes to
abnormal growth
abnormal growth
abnormal growth

signals gets global turned on I-2 & I-3 & come up
and new cell types when stimulate?
I-2 affects growth (remove)

4. 10. 2018

T-test vs Z-test

208
Name (35)

Size of P-value

Y \downarrow No

Z \downarrow

330

Y \rightarrow No

Z \rightarrow No

II

Types of Errors

→ Type-I error — is the rejection of the null hypothesis when it is actually true.

→ Type-II error → fail to reject the null hypothesis that is actually false.

Confidence Interval & Margin of Error

estimate :-

Point \downarrow
Interval

Single Val Point \pm error

800 - 1000
More accurate representation of reality

C.I → What Value the Sample Statistics with take can be known through C.I,

C.I

⇒

CHI SQUARE TEST

- Q) 12% of people are left handed. To verify this theory you took a sample of 75 students, 11 are left handed.

With 5% level of significance, conduct the test.

$$\begin{aligned} \text{H}_0: & \pi = 12\% \\ \text{H}_1: & \pi \neq 12\% \end{aligned}$$

$$\alpha \approx 5\%$$

(3) Chi Sq tested

↳ Chi square Statistics

$$\chi^2_{\text{obs}} = \sum \frac{(O-E)^2}{E}$$

$$\chi^2_{\text{obs}} = \frac{(11-9)^2}{9} + \frac{(64-66)^2}{66} \Rightarrow \frac{2^2}{9} + \frac{2^2}{66} = 0.505$$

$$\Rightarrow \chi^2_{\text{critical}} \text{ for } \alpha = 0.05$$

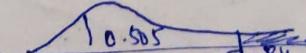
$$\text{2 categorical Variable} \quad \text{dof} = 2 - 1 = 1 \quad (\text{no of categories} - 2)$$

$$\chi^2_{\text{table}} \Rightarrow \text{dof } 1, \alpha = 0.05$$

↓ ↓

$$\chi^2_{\text{critical}} = 3.84$$

Steps
1 tail test



if $\chi^2_{\text{obs}} > \chi^2_{\text{critical}} \rightarrow \text{Reject H}_0$
 $0.505 < 3.84 \rightarrow \text{fail to reject H}_0$

Conclusion:
 12% of people are left handed with 95% confidence.

F-test:

- follow F-distribution
- Comparing Variances
- $F = \frac{S_1^2 / d_1}{S_2^2 / d_2}$ where
- right skewed
- Non-negative
- It depends on DOF

(8) The following data is about the numbers of bulbs produced daily by two workers A and B.

A	B
40	39
30	38
38	41
41	33
38	32
40	39
35	40
	34

- $\alpha = 0.05$
- can we consider based on data that worker B is more stable and efficient?
 - Why not mean can be used for test?
 - mean is same for the both the sample so we will compare variance.

$$\frac{(3-0)^2}{3} = 9$$

$$MSE = \frac{s^2}{n} = \frac{9}{3} = 3 \quad S_{\bar{x}} = \sqrt{\frac{S^2}{n}} = \sqrt{\frac{9}{3}} = \sqrt{3}$$

20.0 > critical value

(3 degrees of freedom) $F_{(2,1)} = 4.25$ at 95% significance

20.0 > 4.25, rejects H₀

∴ H₀ is rejected

∴ P-value < 0.05

∴ reject H₀ with 5%

all terms & final