

Introduction to Statistics

Defn : Stats is the science of collecting, organizing and analyzing data.

Data : "facts or pieces of Information"

Eg : Height of students in classroom

No. of sales in term of revenue of a company

Salary of people in Society.

IQ of students in classroom.

Types Of Statistics

① Descriptive Stats

Defn : It consists of organizing, summarizing and visualizing data.

① Measure of Central Tendency [Mean, Median, Mode]

② Measure of Dispersion [Variance, Standard deviation, Z-score].

③ Different types of Distribution of data

Eg: Histogram, pdf, pmf, Gaussian Distr.,
Log Normal Dist, Exponential, Binomial,
Bernoulli Distr, Poisson

} EDA

Boxplot, Violin plot, Histogram, Distribution Plot

② Inferential Stats

Defn : It consists of using data you have measured to form conclusion.



- 1) Z-test
 - 2) t-test
 - 3) Chi-Square Test
 - 4) Anova Test
- hypothesis
Testing, p-values
Significance
Value.

And [Feature Engineering]

EDA → Exploratory Data Analysis.

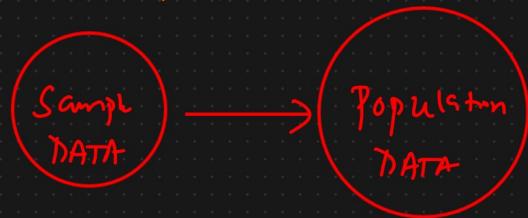
Eg: Let say there are 10 cricket camps in Bangalore. And you have collected the height of cricketers from one of the camp.

Ans) Heights are recorded $[175\text{cm}, 180\text{cm}, 140\text{cm}, 140\text{cm}, 135\text{cm}, 160, 135]$.
↓
Sample

Descriptive Question : [What is the average height of the entire classroom]?
[Distribution of a DATA]
→ [140cm how many std it is away from mean] → Z-score.

Infuential Question

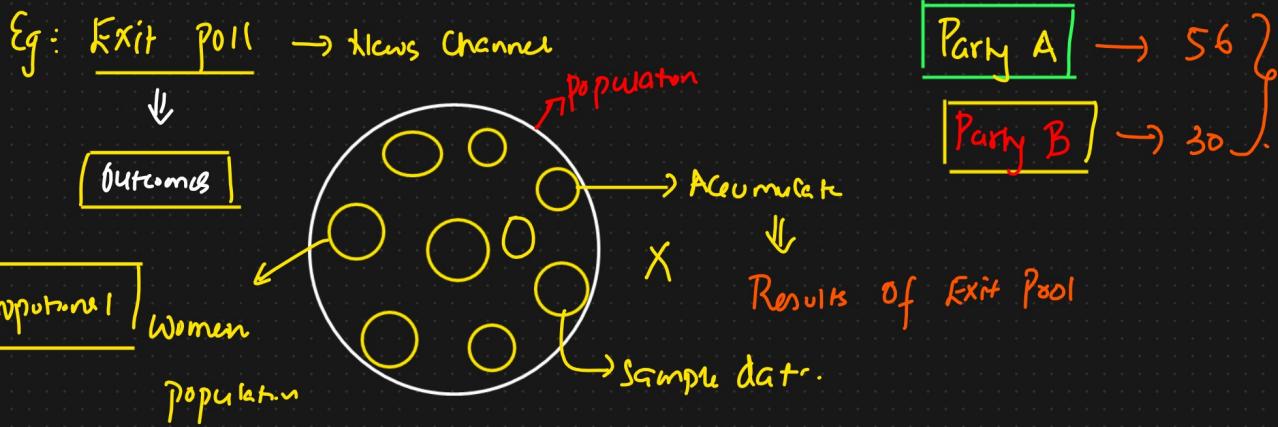
If average
Are the ^{height} of players of camp 1 similar to that of Camp 2



① Population And Sample Data

Population (N) : Population is a group or a superset of data that you are interested in studying.

Sample (n) ! A Sample is a subset of population data.



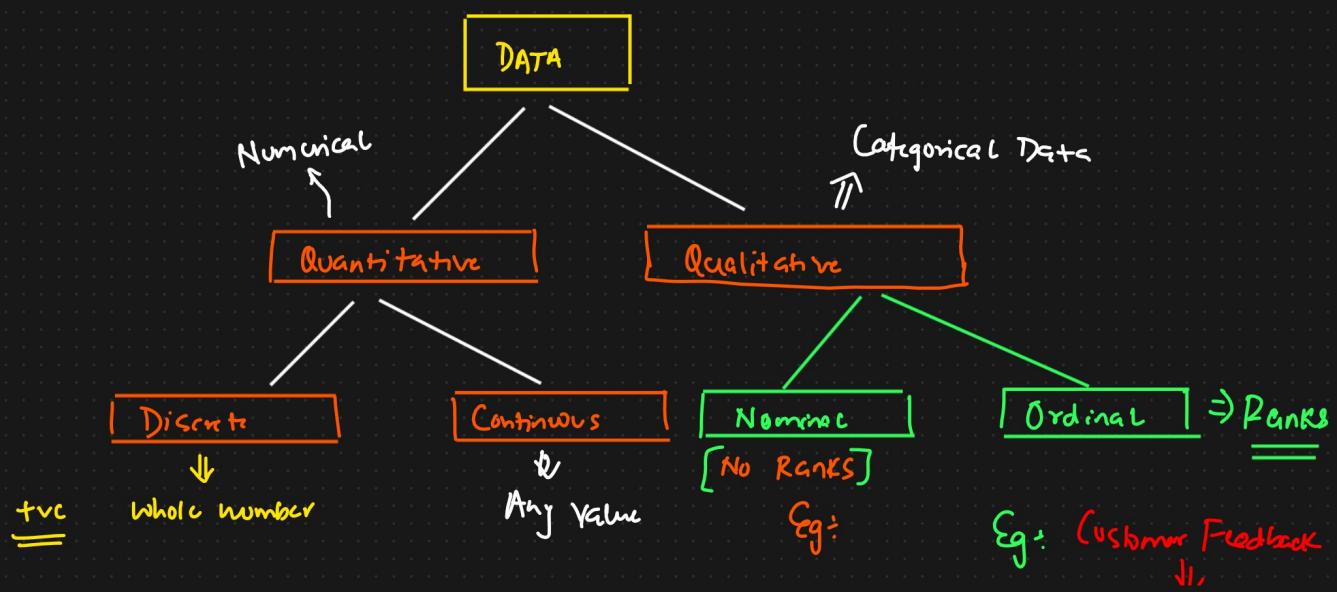
Survey ÷ DATA SCIENCE ⇒ Tech

↳ Farmers

Eg: Types of Sampling Techniques → Assignment

③ Types of Data [House Price Data].

Size	Area	No. of Room	Location	Price
1500	Posh	3	New York	—
—	—	—	Bangalore	—
			New Delhi	



Eg: No. of children in a family.

Eg: Average Rainfall

Eg: Gender,
M, F

{1, 2, 3, 4, 5}.

No. of accounts in banks

House price In Bangalore

Blood group

Weekdays

=====

No. of bikes

Length of River

Country

No. of people working
in a company

Location

Cities.

Days

Pincode.

Example Assignment

- ① Length of Different Rivers In the World?
- ② Favorite Food Based On Gender?
- ③ Marital Status?
- ④ ID measurement?
- ⑤ Weekdays For IT Employees?
- ⑥ Weekdays For KRISH?

Statistics

Agenda

- ① Scale of Measurement ✓
- ② Measure of Central Tendency ✓
- ③ Measure of Dispersion ✓
- ④ Random Variables ✓
- ⑤ Sets ✓
- ⑥ Histogram And Skewness.

① Scale of Measurement

- ① Nominal Scale Data ✓
- ② Ordinal Scale Data ✓
- ③ Interval Scale Data ✓
- ④ Ratio Scale Data

① Nominal Scale Data

- 1) Qualitative/Categorical Data
- 2) Eg: Gender, Color, Labels
- 3) Order or Rank does not matter

Eg: Favorite Colors

Gender

$$\begin{array}{l} \text{Red} \rightarrow 5 \rightarrow 50\% \\ \text{Yellow} \rightarrow 3 \rightarrow 30\% \\ \text{Orange} \rightarrow 2 \rightarrow 20\% \\ \hline 10 \end{array}$$

② Ordinal Scale Data

- 1) Rank is important
- 2) Order matters
- 3) Difference cannot be measured

Eg: Race [100m]

$$\begin{array}{l} \rightarrow 1^{\text{st}} \\ \rightarrow 2^{\text{nd}} \\ \rightarrow 3^{\text{rd}} \end{array}$$

$\left\{ \begin{array}{l} 1^{\text{st}} \\ 2^{\text{nd}} \\ 3^{\text{rd}} \end{array} \right\}$

| IQ |

③ Interval Scale Data

Eg: Temperature Variable

↓
[2:1]

① The order matters ✓ [Rank].

30°F
60°F

$$60:30 = 2:1$$

② Difference can be measured ✓

90°F

[OF] X

°C

③ The Ratio Cannot be measured ✓

120°F

④ No "0" Starting point

④ Ratio Scale Data

Eg: Students marks in a class

0, 90, 60, 30, 75, 40, 50
= [3:2] [3:1]

① The order matter

② Differences are measurable (Ratio)

③ Contains a "0" Starting point

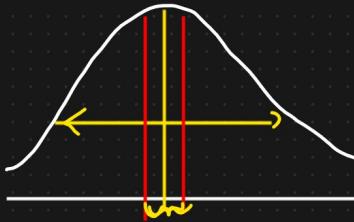
② Measure of Central Tendency [EDA] → Exploratory Data Analysis

① Mean or Average

② Median

③ Mode

}



① Mean

Population (N)

$$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$

$$\text{Population mean} (\mu) = \sum_{i=1}^N \frac{x_i}{N}$$

$$= \left[1 + 1 + 2 + 2 + 3 + 3 + 4 + 5 + 5 + 6 \right]$$

$$= \frac{32}{10} = 3.2$$

Sample (n)

$$X = \{2, 3, 1, 4, 8, 9, 5\}$$

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

$$= \frac{32}{7} = 4.7$$

② Median

$$X = \{4, 5, 2, 3, 2, 1\}.$$

① Sort the Random Variable

- ② No. of Elements
- If $n = \text{even} \Rightarrow$ Middle 2 elements \Rightarrow Average
 - If $n = \text{odd} \Rightarrow$ Central Element

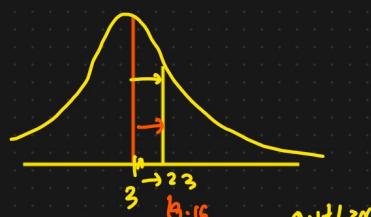
$$X = \{1, 2, \boxed{2, 3}, 4, 5\}.$$

$$\downarrow \\ \frac{2+3}{2} = 2.5.$$

$$\{1, 2, 3, 4, 5\}$$

$$\boxed{\text{Mean} = 3}$$

$$\boxed{\text{Median} = 3}$$



$$\{1, 2, \boxed{3}, 4, 5, 100\}.$$

$$\downarrow \\ \boxed{\text{Mean} = 19.16}$$

$$\downarrow \\ \boxed{\text{Median} = 3.5}$$

$$\frac{1+2+3+4+5+100}{6} = \frac{115}{6} = 19.6$$

TRANSACTION FRAUDULENT {Outliers}.

$$X = \{2, 8, 4, 5, 1, 7, 9, 120, 130\}.$$

Mean And Median

Outlier

$$\{1, 2, 4, 5, \boxed{7, 8, 9, \boxed{120, 130}}\}$$

Median

$$\text{Mean} = \frac{1+2+4+5+7+8+9+120+130}{9} = 31.77$$

$$\boxed{\text{Median} = 7}$$

$$\frac{37}{7} = 5.1$$

Outliers of Median \rightarrow calculate central tendency}.

$$X = \left\{ \boxed{-5, -10} \right\} \cup \{1, 2, 3, 4, 5\}. \Rightarrow \{ -10, -5, \boxed{1}, \boxed{2}, 3, 4, 5 \}$$

$$\underline{\text{Mean}} = 0 \quad \underline{\text{Mean}} = 3$$

$$\underline{\text{Median}}$$

③ Mode : Frequency of maximum occurring element.

$$X = \left\{ \underbrace{2, 1, 1, 1}_{\Downarrow}, 4, 5, 7, 8, 9, 9, 10 \right\}.$$

$$\underline{\text{Mode}} = 1.$$



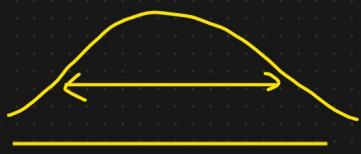
EDA And FE

→ Missing Value or Null Values

Age	Weight	Salary	Gender	Degree
24	70	40K	M-	BE
25	80	70K	F	-
27	95	45K	F	-
24	-	50K	M-	PhD
32	-	60K	{Mode}	B.E
<u>Mean</u>	60	-	-	Marks
	65	55K	-	BSC
	40	72	M-	B.E

$$\underline{\text{Median}} = 150$$

② Measure of Dispersion [Spread of the data].



① Variance

② Standard Deviation

① Variance

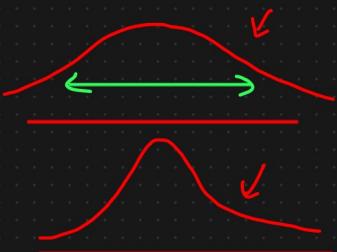
Population Variance

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \quad \{ \text{Spread} \}$$

$x_i \rightarrow$ DATA Points

$\mu \rightarrow$ population mean

$N \rightarrow$ Population Size



Sample Variance

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

$x_i \rightarrow$ DATA POINTS

$\bar{x} \rightarrow$ Sample Mean

$n \rightarrow$ Sample Size

Assg

Why we divide sample variance by $n-1$?

↑ Bessel's Correction.

Ans) The sample variance is divide by $n-1$ so that we can

Create an unbiased estimator of the population variance

Population Std

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

Sample Std

$$s = \sqrt{\frac{\text{Sample Variance}}{n-1}}$$

How many std x_i is away from Mean

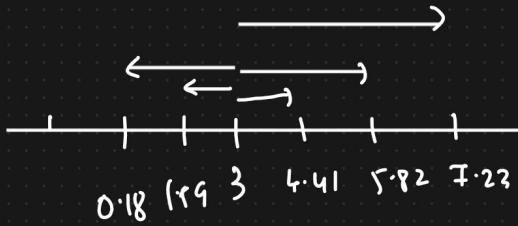
Eg: $\{1, 2, 3, 4, 5\}$

$$\mu = 3$$

$$\text{Population Variance} = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}$$

$$= \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2$$

$$\sigma = \sqrt{2} = 1.41$$



$$\begin{array}{r}
 3'00 \\
 1.41 \\
 \hline
 1.59 \\
 1.41 \\
 \hline
 0.18 \\
 \hline
 5.82 \\
 1.41 \\
 \hline
 7.23
 \end{array}$$

Variables (x, y)

4) Random Variables

$$x+5=7$$

$$\boxed{x=2}$$

$$\boxed{x=2}$$

$$\boxed{y=6}$$

$$8=y+x$$

Exponent
↑

Random Variable is a process of mapping the output of a random process or experiment to a number.

Eg: Tossing a Coin

$$X = \begin{cases} 0 & \text{if Head} \\ 1 & \text{if Tail} \end{cases}$$

Rolling a dice 6 times

$$Y = \{2, 1, 4, 2, 3, 5\}$$

$Y = \{ \text{Sum of Rolling a dice } 7 \text{ times} \}$.

$$\Pr(Y > 15) \quad \Pr(Y < 10)$$

⑤ Sets

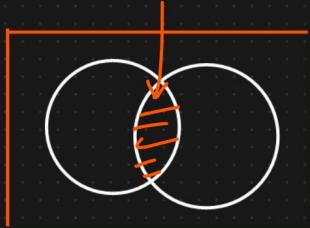
$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7\}$$

① Intersection

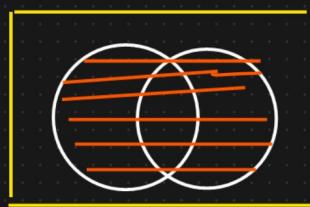
$$A \cap B = \{3, 4, 5, 6, 7\}$$

$A \cap B$.



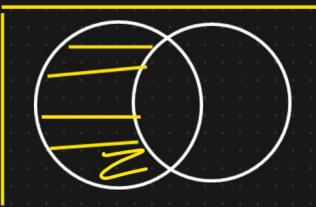
② Union

$$A \cup B = \{1, 2, 3, 4, 5, 6, 7, 8\}$$



③ Difference

$$A - B = \{1, 2, 8\}$$



④ Subset

$$A \rightarrow B \Rightarrow \text{False}$$

$$B \rightarrow A \Rightarrow \text{TRUE}$$

⑤ Superset

$A \rightarrow B \Rightarrow \text{TRUE}$

$B \rightarrow A \Rightarrow \text{False}$.

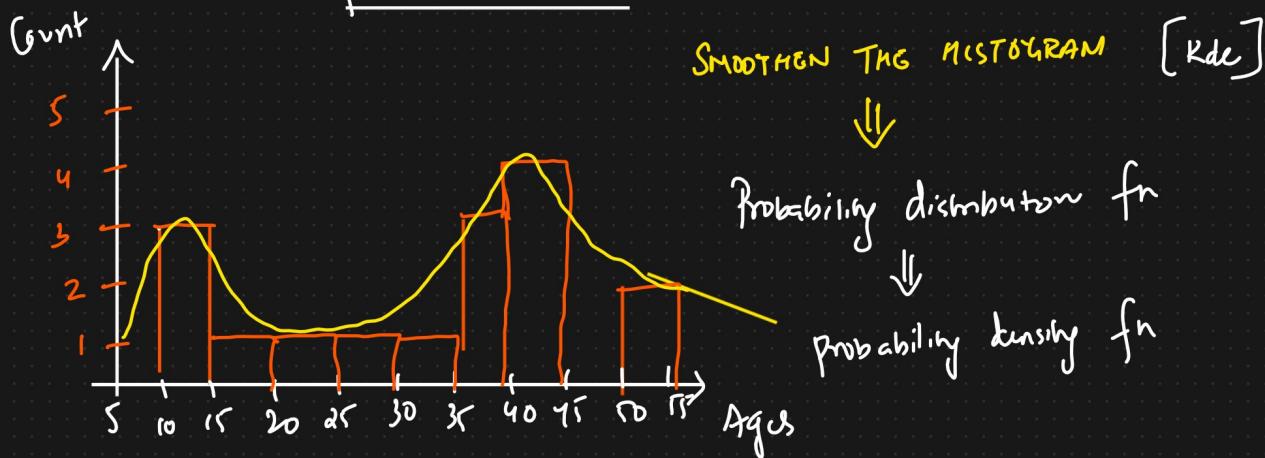
Histograms And Skewness \rightarrow [frequency]. \Rightarrow Distribution

$$\text{Age} = \{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51\}.$$

Bins, Bin Size

$$\frac{50}{5} = \boxed{10} \Rightarrow \text{No. of bins}$$

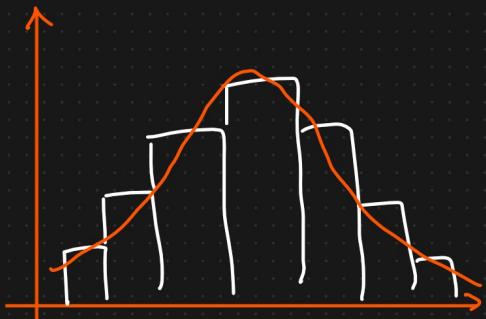
$$\boxed{\text{Bins size} = 5}$$

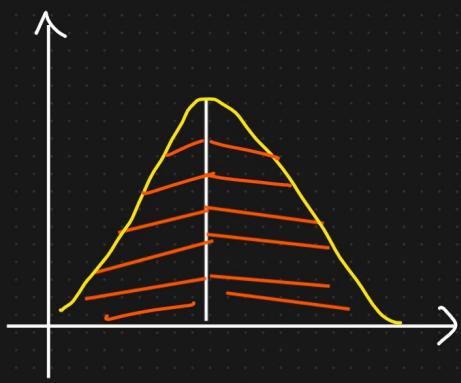


Skewness

Normal/Gaussian Distribution

\Rightarrow Symmetrical Distribution





① No Skewness

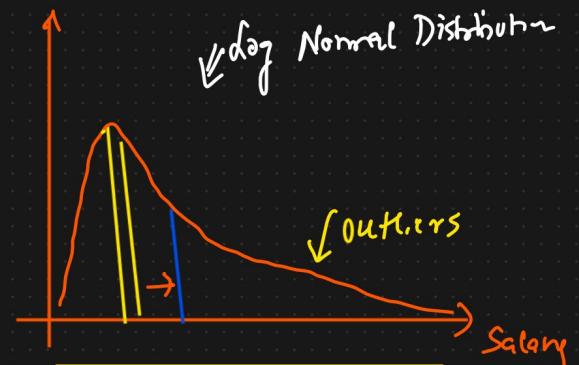
Mean, Median, Mode

$$\boxed{\text{Mean} = \text{Median} = \text{Mode}}$$

② Right Skewed

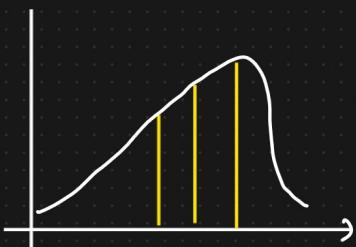


\Rightarrow Positive Skewed \Rightarrow



$$100, 200, 300, 400, 1000000 \rightarrow \boxed{\text{mean} > \text{median} > \text{mod}}$$

③ Left Skewed



\Rightarrow Negative Skewed

Relationship : $\text{mean} \leq \text{median} \leq \text{mode}$

Statistics

Agenda

- ① Sampling Techniques ✓
- ② Covariance And Correlation
- ③ Probability distribution function
- ④ Probability density function
- ⑤ Probability Mass function
- ⑥ Cumulative Density Function
- ⑦ Types of Distribution

① Sampling Techniques

① Random Sampling : Simple Random Sampling gives each member of the population an equal chance of being chosen for the sample.

Eg: Vaccination Test \rightarrow 100 \rightarrow Randomly Select Couple

Accidental Test \rightarrow 1000 \rightarrow A \rightarrow 2 vehicle

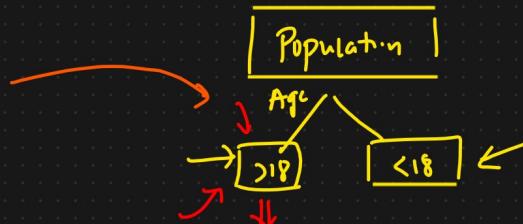
Exit Poll

Average IQ of the school \rightarrow Select 10 people \Rightarrow

② Stratified Sampling [Stratified \rightarrow layers]

Stratified Sampling Involves dividing the population into sub population that may differ in important ways

Eg: Exit Poll



Voting \rightarrow Random Sampling.



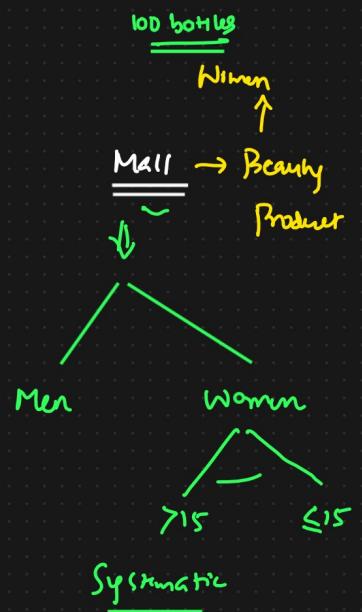
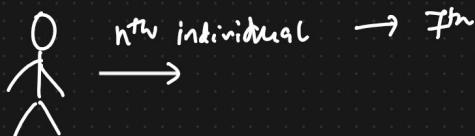
Hale
Funk

③ Systematic Sampling : Systematic Sampling is a statistical method involving the selection of elements from an ordered sampling frame.

Eg: Airport

A Credit Card Company

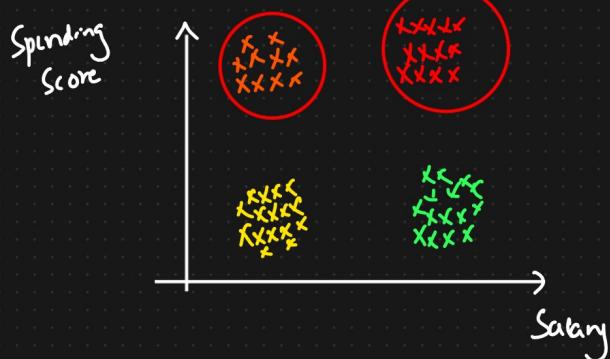
Mall → Fill up a form for
a cause.



④ Convenience Sampling Assignment

⑤ Purposive Sampling → Judgemental Sampling is a method where researchers decides which members of the target population will be sampled.

⑥ Cluster Sampling



Mail iphone 28

→ 100 ↗
→ 20r.

② Covariance And Correlation

X Y

2 3

4 5

6 7

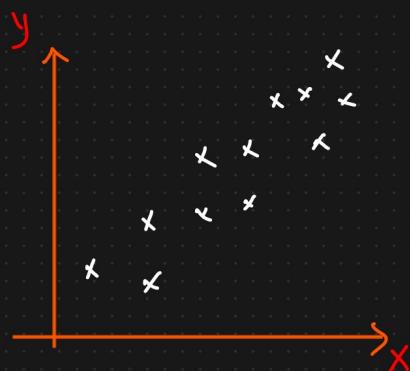
8 9

$X \uparrow$	$y \uparrow$
$X \downarrow$	$y \downarrow$

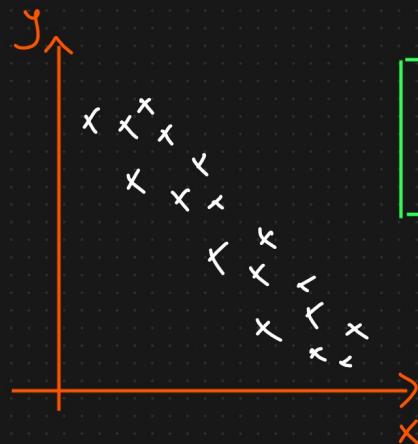
Quantify
[Relationship between X And y].



Covariance And Correlat.-r.



$X \uparrow$	$y \uparrow$
$X \downarrow$	$y \downarrow$



$X \uparrow$	$y \downarrow$
$X \downarrow$	$y \uparrow$

Covariance

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

x_i → Data point of x

\bar{x} → Sample mean of x

y_i → Data points of y

\bar{y} → Sample mean of y .

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

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

$$\text{Var}(x) = \text{Cov}(x, x) \Rightarrow \underline{\underline{\text{Sprod}}}$$

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

$X \uparrow$	$y \uparrow$
$X \downarrow$	$y \downarrow$

+ve Covariance

$X \uparrow$	$y \downarrow$
$X \downarrow$	$y \uparrow$

-ve Covariance

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

$\rightarrow 2 \quad 3$

$\rightarrow 4 \quad 5$

$\underline{6} \quad \underline{7}$

$\bar{x}=4 \quad \bar{y}=5$

$$= \left[\frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{2} \right]$$

$$= \frac{4+0+4}{2} = 4 \text{ //.}$$



$X \quad Y$

$8 \quad 6$

$7 \quad 5$

$6 \quad 4$

$3 \quad 1$

$2 \quad 0$

$$\bar{x} = 5.2 \quad \bar{y} = 3.2$$

$X \quad Y$

$8 \quad 6$

$9 \quad 5$

$10 \quad 4$

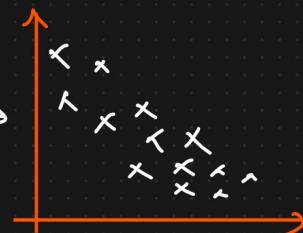
$12 \quad 2$

$11 \quad 1$

$$\bar{x} = 10 \quad \bar{y} = 3.6$$

-ve Covariance

$$[-3]$$



$$\boxed{6.7}$$

Advantages

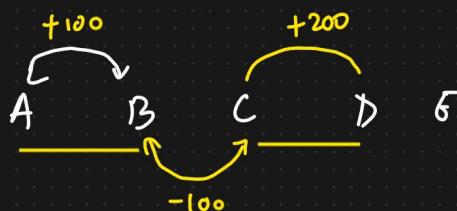
① Relationship between X and Y

+ve or -ve value

$A \leftrightarrow B$ $+ \infty \quad - \infty$

Disadvantages

① Covariance does not have a specific limit value



5

② Pearson Correlation Coefficient [-1 to 1].

$$\rho_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y} \Rightarrow -1 \text{ to } 1$$

① The more the value towards +1, the more +ve correlated x, y is.

② The more the value towards -1, the more -ve " x, y is

$$\begin{array}{cc} x & y \\ 8 & 6 \end{array}$$

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

$$\begin{array}{cc} 9 & 5 \end{array}$$

$$\text{Cov}(x,y) = -3$$

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

$$\begin{array}{cc} 10 & 4 \end{array}$$

$$\begin{array}{cc} 12 & 2 \end{array}$$

$$\rho_{x,y} = \frac{-3}{1.58 \times 2.073} = \frac{-3}{3.27}$$

$$\begin{array}{cc} 11 & 1 \\ \hline \bar{x} = 10 & \bar{y} = 3.6 \end{array}$$

$$\left[\frac{4+1+0+4+1}{4} \right] = \frac{10}{4} = 2.5$$

$$\sigma = \sqrt{2.5}$$

$$\sigma_x = 1.58 \quad \sigma_y = 2.073 \quad = -0.917 \Rightarrow \text{Negative Correlated} \Rightarrow 91.7\%$$

$$x \uparrow \quad y \downarrow \Rightarrow 91.7\%$$

$$1 \text{ unit } x \uparrow \quad y \downarrow \quad 0.917$$

③ Spearman Rank Correlation $[-1 \text{ to } 1]$.

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

X	Y	R(X)	R(Y)
1	2	2	2
3	4	3	3
5	6	4	4
7	8	5	6
0	7	1	5
8	1	6	1

Ascending Order.

{ Non linear Relationship }

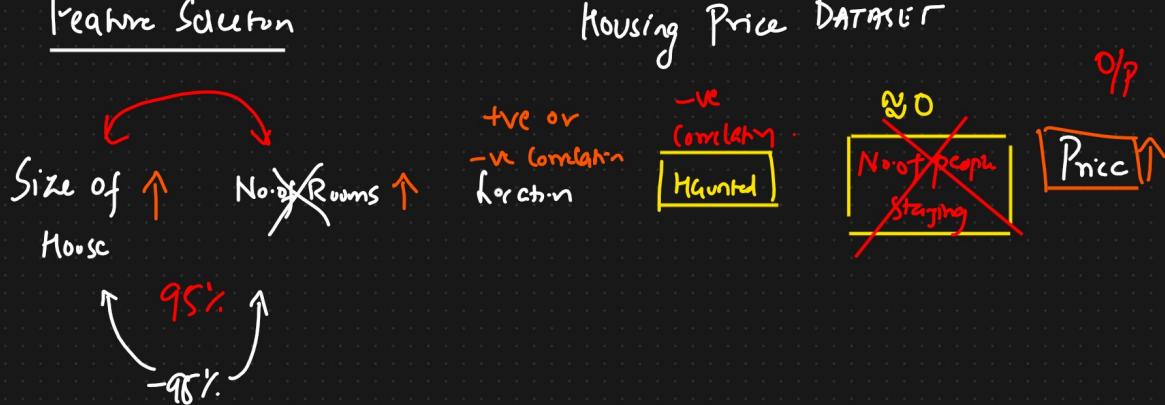
Spearman Rank Correlation



Non linear Relationship

DATA SCIENCE

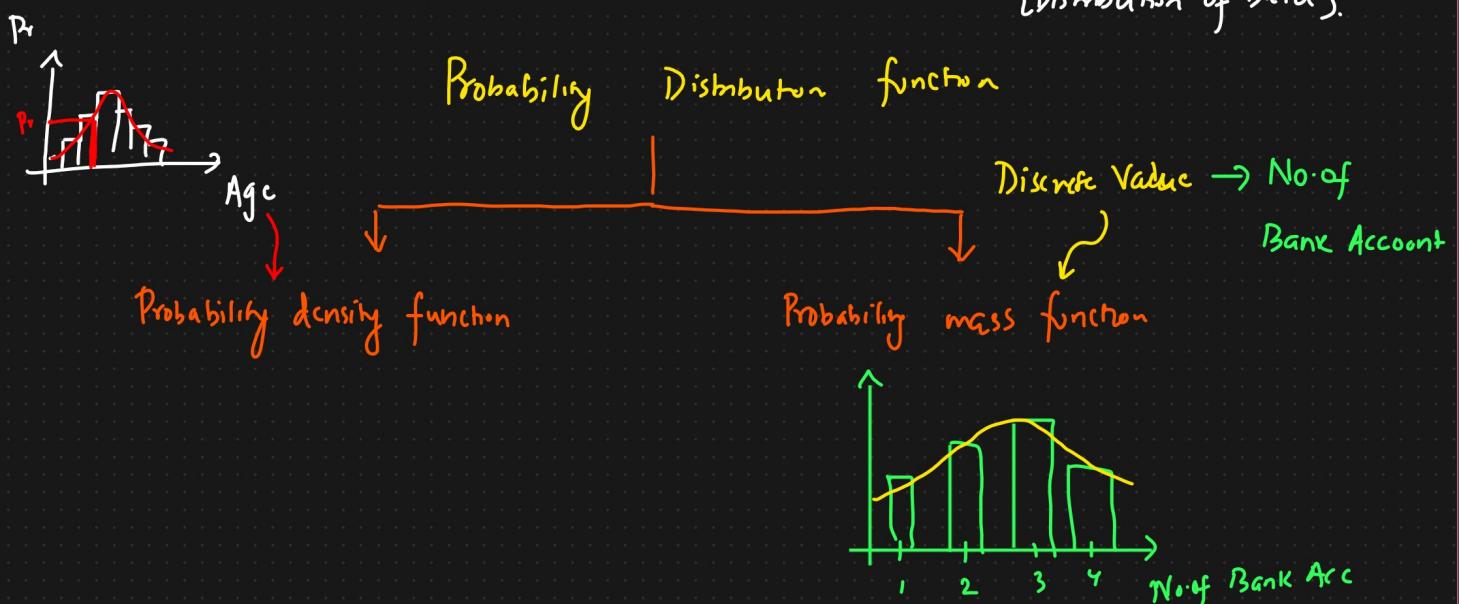
Feature Selection



Statistics

- ① Probability Distribution function
- ② Probability Density function (pdf)
- ③ Probability mass function (pmf)
- ④ Different types of distribution
- ⑤ CDF {Cumulative distribution fn}.

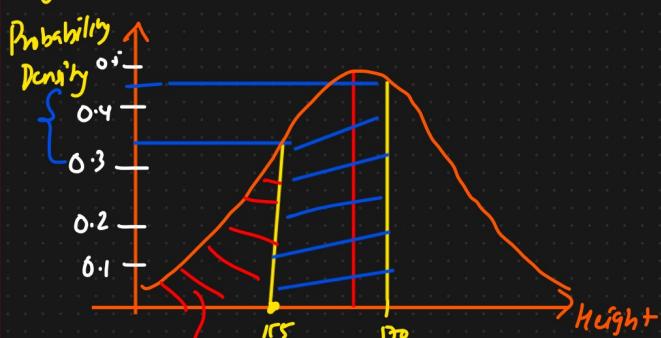
$X = \{\text{Continuous}\}$ or Discrete Values
 {Distribution of Data}.



① Probability Density function

① Continuous Random Variable

Eg: Height of Students in classroom



Formula:
 ↑
 ↓

$$\boxed{\Pr(H \leq 155)} \Rightarrow \underline{\text{Formula}}$$

$$\Pr(H > 155 \text{ and } H \leq 170) =$$

↓ 165

Area under this curve

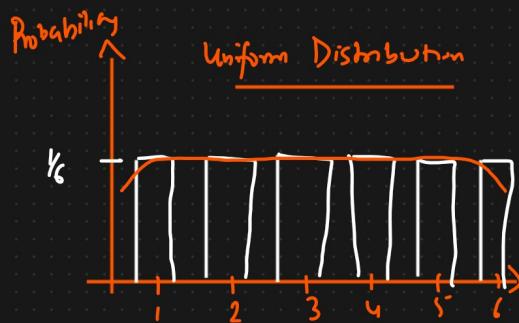
② Probability Mass Function (pmf)

① Discrete Random Variable

Eg: Rolling a Dice $\{1, 2, 3, 4, 5, 6\}$

$$\Pr(1) = \frac{1}{6} \quad \Pr(3) = \frac{1}{6}$$

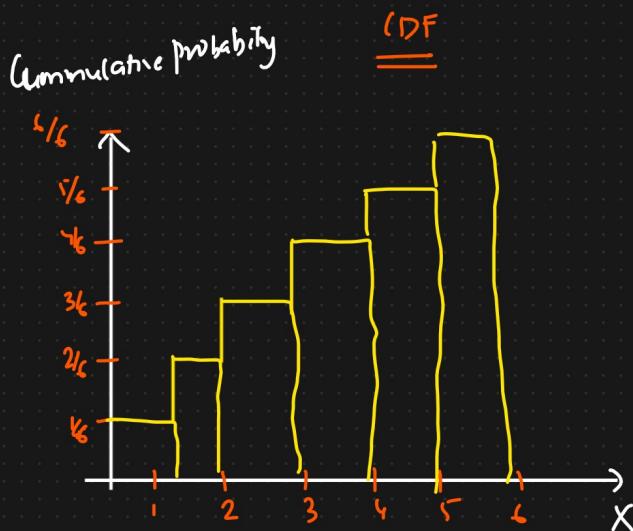
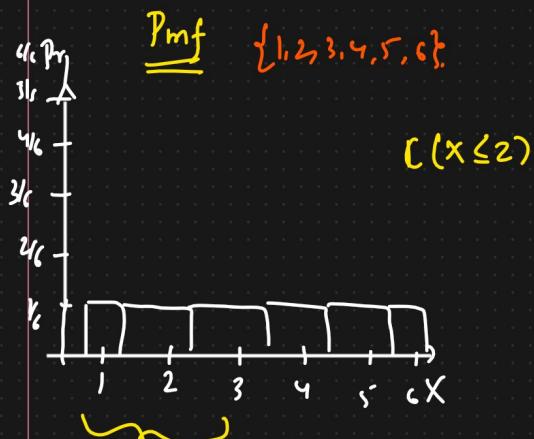
$$\Pr(2) = \frac{1}{6} \quad \Pr(4) = \frac{1}{6}$$



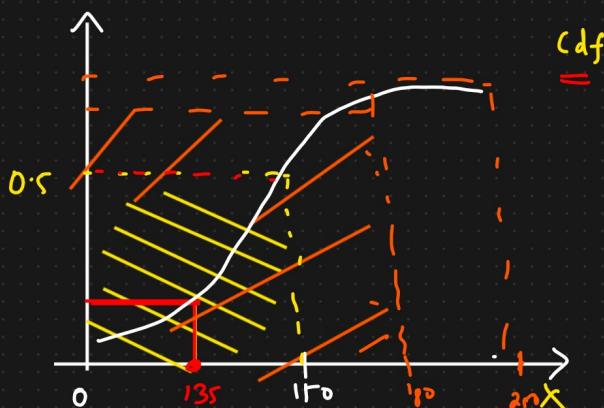
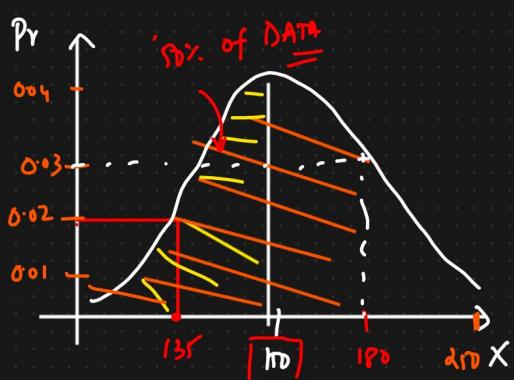
$$\begin{aligned}\Pr(X \leq 4) &= \Pr(X=1) + \Pr(X=2) + \Pr(X=3) \\ &\quad + \Pr(X=4) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{6} = \frac{2}{3}.\end{aligned}$$

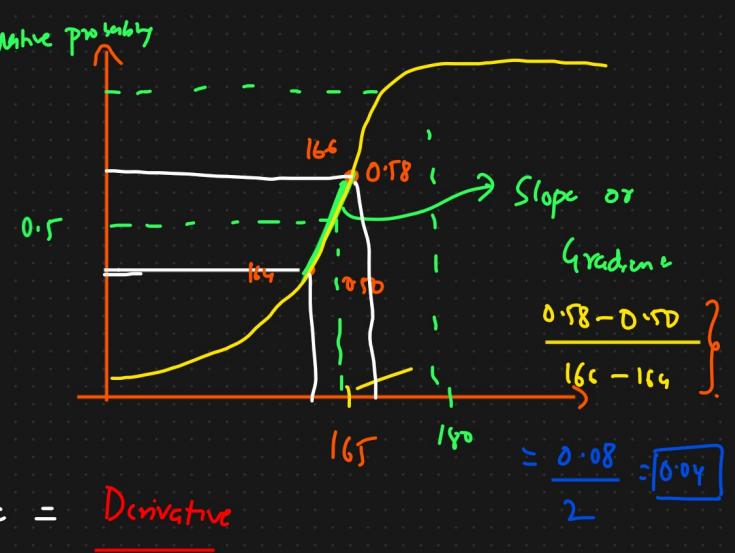
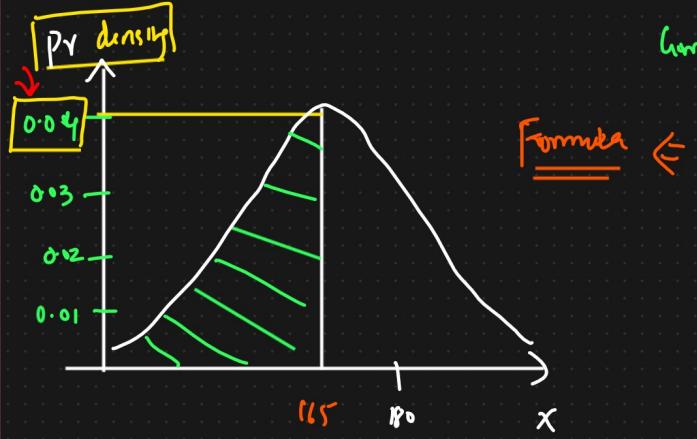
③ Cumulation Distribution Function

Eg: Rolling a dice



② pdf and cdf





$$\text{Slope} = \frac{\text{Derivative}}{ }$$

Probability Density of pdf \Rightarrow Gradient or derivative of cdf

pdf is the derivative of cdf.

cdf is the integration of pdf.

Types of Probability Distribution

- ① Normal/Gaussian Distribution (pdf) ✓
- ② Bernoulli Distribution (pmf) \rightarrow Outcomes are binary only. ✓
- ③ Uniform Distribution (pmf)
- ④ Poisson Distribution (pmf) ✓
- ⑤ Binomial Distribution (pmf) ✓
- ⑥ Log Normal Distribution (Pmf)

① Bernoulli Distribution

① Discrete Random Variable {pmf}

② Outcomes are Binary

Eg: Tossing a coin {H, T}.

$$\Pr(H) = 0.5 = P_H$$

$$\Pr(T) = 0.5 = 1 - p = q_H$$

$P \Rightarrow$ probability of one outcome

$q \Rightarrow$ probability of other

outcome

Eg: Whether the person will pass/Fail

$$\Pr(\text{Pass}) = 0.7 = p$$

$$\Pr(\text{Fail}) = 1 - p = 0.3 = q$$

$$\begin{array}{c} \uparrow \\ p(T) = 0.8 \Rightarrow p \\ q = 1 - 0.8 = 0.2 \end{array}$$

Parameters

PMF

$$\text{pmf} \quad \left\{ \begin{array}{ll} q = 1 - p & \text{if } k=0 \\ p & \text{if } k=1 \end{array} \right\}$$

$$0 \leq p \leq 1$$

$$q = 1 - p$$

$K = \{0, 1\} \Rightarrow$ Outcomes



Formula

Heads and Tail

$$K = 0, 1$$

$$\boxed{\text{PMF} = p^K * (1-p)^{1-K}}$$

$$K \in \{0, 1\}$$

$$P \quad q \quad \text{if } K=1 \quad \text{if } K=0$$

$$\Pr(H) \quad \Pr(T)$$

$$\text{PMF} = P_H$$

$$\begin{aligned} \text{pmf} &= p^0 * (1-p)^1 \\ &= (1-p) = q_H \end{aligned}$$

② Mean of Bernoulli Distribution

$$E(K) = \sum_{k=0}^K k \cdot p(k)$$

$$= [0 \cdot 0.4 + 1 \cdot 0.6]$$

$= 0.6 = P_{H.}$

$$\Pr(K=1) = 0.6 \Rightarrow P$$

$$\Pr(K=0) = 0.4 \Rightarrow 1-P$$

$\{1, 2, 3, 4, 5, 6\}$

④ Binomial Distribution

Note → Refer Wikipedia

- ① Discrete Random Variable.
- ② Every Experiment outcome is Binary.
- ③ This experiment is performed for n trials

Eg: Tossing a coin 10 times

Notation: $B(n, p)$

$$\Pr(H) = 0.5 = p$$

$$\Pr(T) = 0.5 = q = 1-p.$$

Parameters: $n \in \{0, 1, 2, 3, \dots\} \Rightarrow$ no. of trials

$p \in \{0, 1\} \rightarrow$ success probability for each trial

$$q = 1-p.$$

Support = $K \in \{0, 1, 2, \dots, n\} \Rightarrow$ Number of Success.

PMF

$$Pr(K, n, p) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

for $k=0, 1, 2, \dots, n$ where

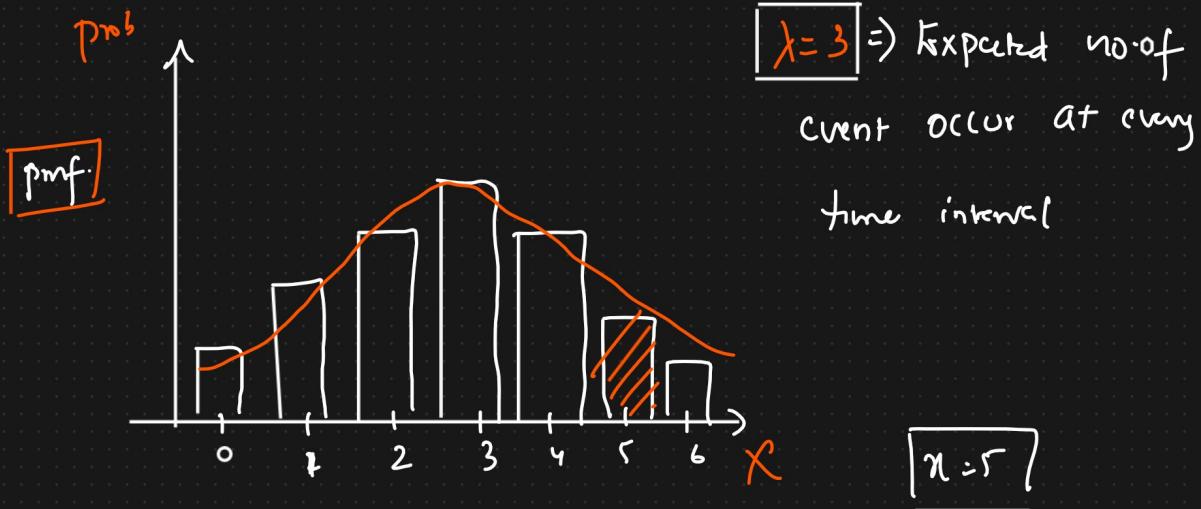
Poisson Distribution

(1) Discrete Random Variable (pmf)

(2) Describe the number of event occurring in a fixed time interval //.

Eg: No. of people visiting hospital every hour

No. of people visiting banks at 11am \equiv  Threshold



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

$$\boxed{\lambda}$$

$\lambda=3 \Rightarrow$ Let us consider $\lambda=3$

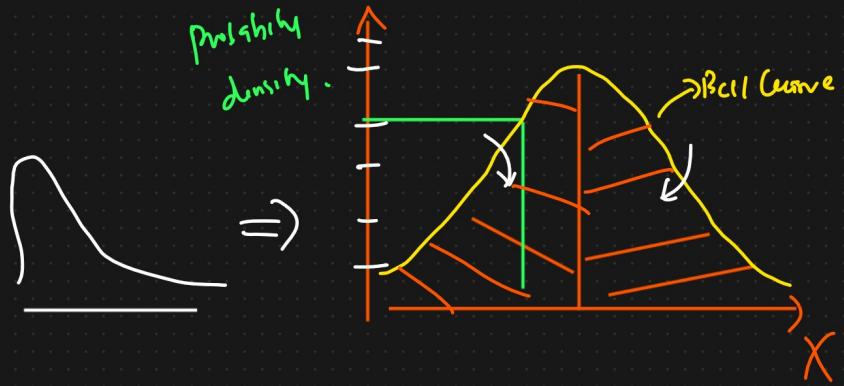
$$= \frac{e^{-3} 3^5}{5!} = 0.101 = \underline{\underline{10.1\%}}$$

4pm and 5pm.

$$\Pr(X=4 \text{ and } X=5) = \Pr(X=4) + \Pr(X=5)$$

* Normal / Gaussian Distribution (pdf)

(*) Continuous Random Variable.



Age, Height, weight.

$X \stackrel{\text{def}}{=} \dots \dots \dots \dots \dots$
↓
D.

[IRIS]

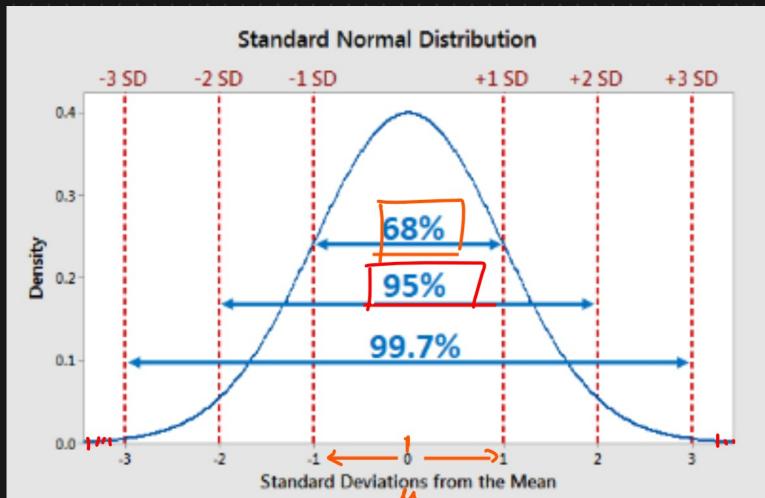
$$\text{Notation} : N(\mu, \sigma^2)$$

Parameters : $\mu \in \mathbb{R}$
 $\sigma^2 \in \mathbb{R} > 0$ = variance
 $\mu \in \mathbb{R}$

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

$$\text{Mean} = \mu \quad \text{Variance} = \sigma^2 \quad \text{Std} = \sqrt{\text{Variance}}$$

Empirical Rule of Normal Distr



$X = \{ \dots, \underline{\quad}, \dots \}$
 \downarrow

Normal
Distribu^r

68% of the entire
distribution will
fall in 1 st
std.

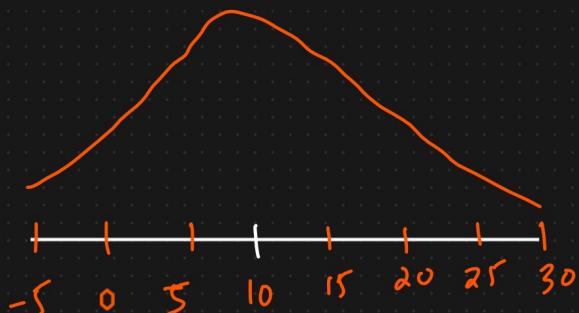
$$\mu - \sigma \leq x \leq \mu + \sigma \approx 68\%$$

$$\mu - 2\sigma \leq x \leq \mu + 2\sigma \approx 95\%$$

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

$$\mu = 10 \quad \sigma = 5 \Rightarrow \text{Normal Distribution}$$

$$\boxed{25}$$



$$f_1 = \quad f_2 = \quad f_3 =$$

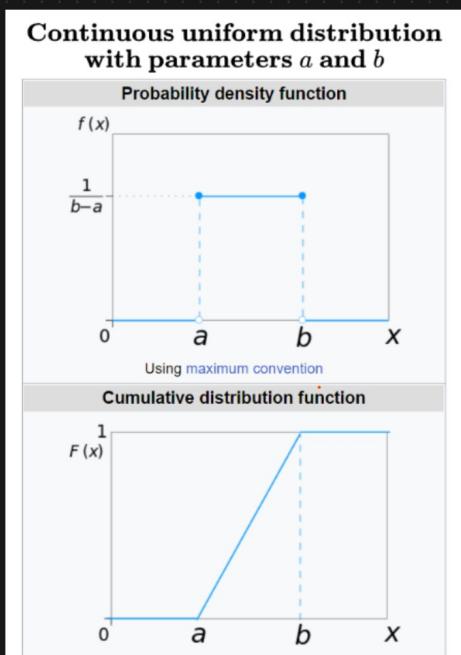
Statistics

Agenda

- ① Uniform Distribution
- ② Standard Normal Distribution And Z-score
- ③ Central limit Theorem
- ④ Log Normal Distribution

① Uniform Distribution

- ① Continuous Uniform Distribution (pdf)
- ② Discrete Uniform Distribution (pmf)
- ③ Continuous Uniform Distribution [Continuous Random Variable]



Notation = $U(a,b)$

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

$$\text{Pdf} = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$$

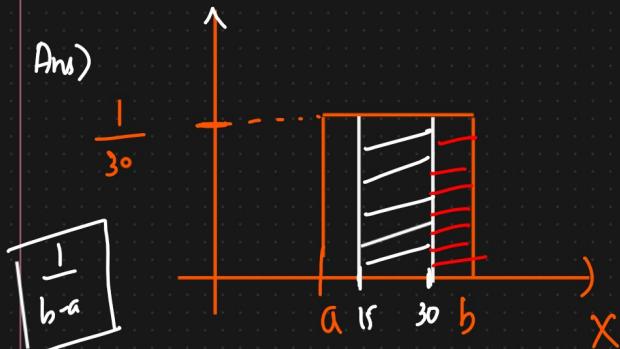
$$\text{cdf} = \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a,b] \\ 1 & \text{for } x > b \end{cases}$$

$$\text{Mean} = \frac{a+b}{2} \quad \text{Median} = \frac{a+b}{2}$$

Eg: The number of candies sold daily at a shop is uniformly distributed with a maximum of 40 and a minimum of 10.

(1) Probability of daily sales to fall between 15 and 30?

Ans)



$$a = 10 \quad b = 40$$

$$x_1 = 15$$

$$x_2 = 30$$

$$\begin{aligned} P(15 \leq x \leq 30) &= (x_2 - x_1) * \frac{1}{b-a} \\ &= 15 * \frac{1}{30-10} = 0.5 \end{aligned}$$

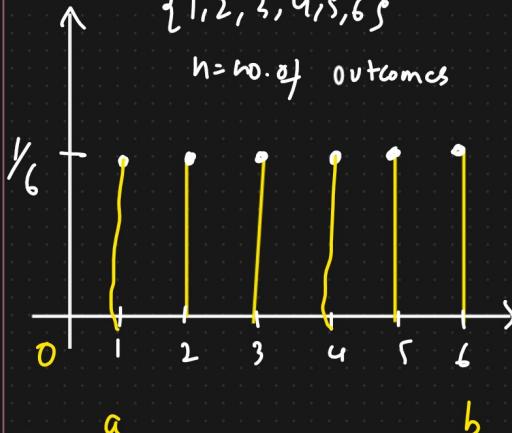
$$P(X > 30) = (40-30) * \frac{1}{30} = 10 * \frac{1}{30} = \frac{1}{3} = 0.33 = 33\%$$

(2) Discrete Uniform Distribution (pmf)

Eg: Rolling a dice

$$\{1, 2, 3, 4, 5, 6\}$$

$n = \text{no. of outcomes}$



$$n = 6$$

$$n = b - a + 1$$

$$= 6 - 1 + 1$$

$$n = 6$$

==

Notation $U(a, b)$

Parameter a, b $\boxed{b > a}$

$$\text{PMF} = \frac{1}{n}$$

$$\text{Mean} = \frac{a+b}{2}$$

$$\text{Median} = \frac{a+b}{2}$$

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

④ Standard Normal Distribution And Z-Score

Z-Score

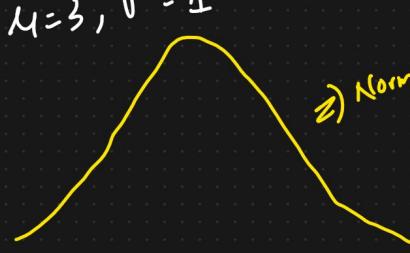
$$X = \{1, 2, 3, 4, 5\}$$

⇒ Normally Distributed

$$\mu = 3$$

$$\sigma = 1.414 \approx 1$$

$$\mu = 3, \sigma = 1$$

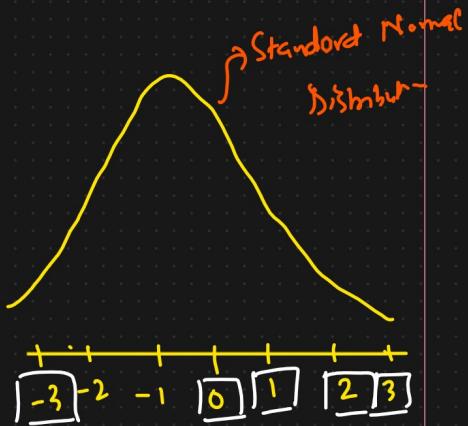


⇒ Normal Distribution

$$\Rightarrow Z\text{-score} = \frac{X_i - \mu}{\sigma}$$

$$X = \{1, 2, 3, 4, 5\}$$

$$\mu = 0 \quad \& \quad \sigma = 1$$



$$Z\text{-score} = \frac{1 - 3}{1} = -2$$

$$\frac{3 - 3}{1} = 0$$

$$\frac{0 - 3}{1} = -3$$

$$\frac{2 - 3}{1} = -1$$

$$\frac{4 - 3}{1} = 1$$

$$\frac{5 - 3}{1} = 2$$

↓
(year)
Age

↓
(kg)
Weight

↓
(cm)
Height

↓
(IN2)
Salary

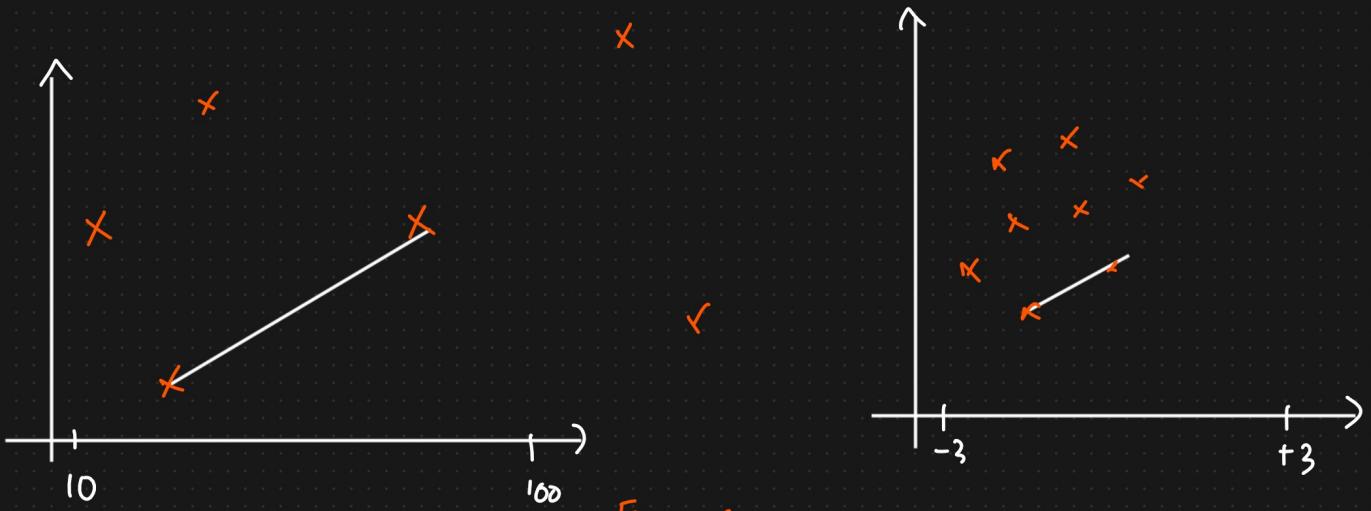
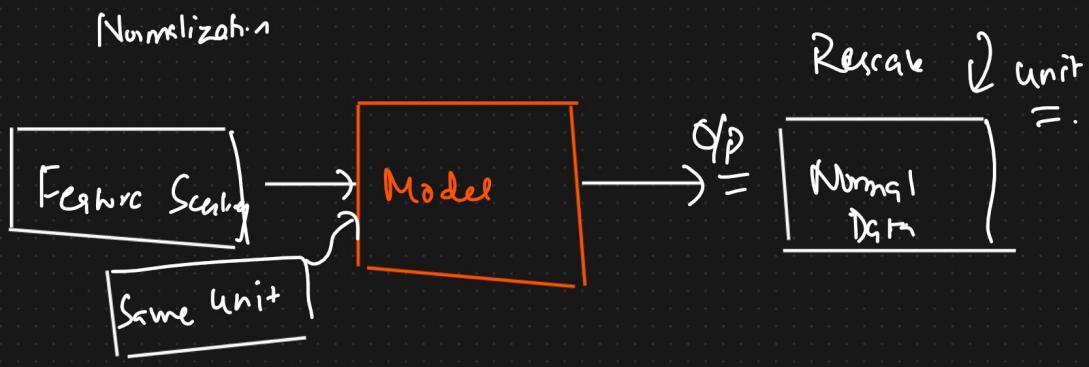
$$\mu = 0 \quad \sigma = 1$$

-3	24	70	-
+3	32	80	-
to	36	75	-
38	60	-	-
20	30	-	-

⇒ Model \Rightarrow Mathematical
Equation

Feature Scaling

$$Z\text{-Score} :=$$



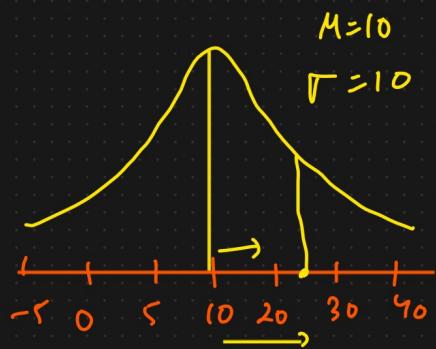
$(-1, 1)$ $\leftarrow (0, 1) \Rightarrow \text{Deep learning}$
 Normalization And Standardization [Assignment].

TRANSFORMATION
 \downarrow

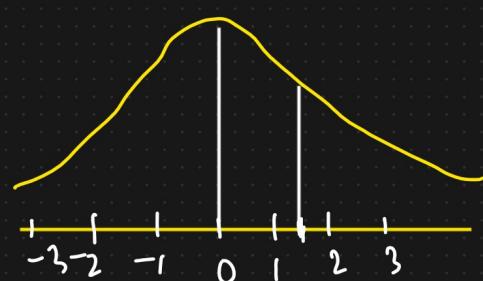
Assignment.

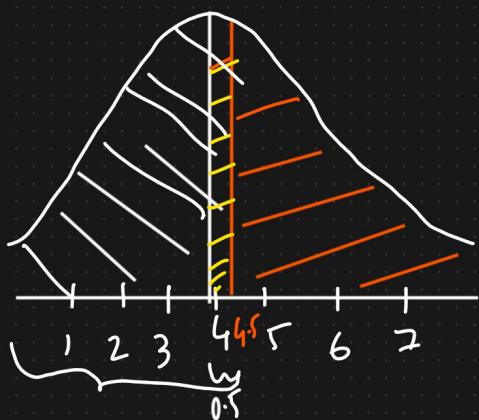
$$\begin{array}{c} \uparrow \downarrow \uparrow \\ \boxed{Z\text{-Score}} \end{array} \quad \mu=0, \sigma=1$$

Feature TRANSFORMATION



$$Z\text{-Score} = \frac{25 - 10}{10} = \frac{15}{10} = 1.5$$





① How many standard deviation 4.5 is away from mean? $\Rightarrow 0.5$

② What percentage of data is falling above 4.5 ?

Z table \Rightarrow Area Under the Curve

$$Z\text{-score} = \frac{4.5 - 4}{1} = 0.5$$

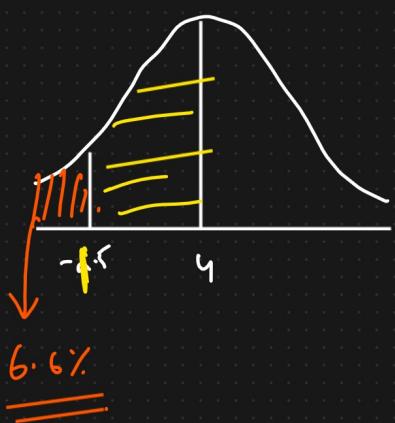
$=$

0.6915

$$\text{Area Under the Curve } (4.5) = 1 - 0.6915 = 0.3085 = 30.85\%$$

① Percentage of data falling below 2.5 ?

$$Z\text{-score} = \frac{2.5 - 4}{1} = -1.5$$



Z table $= 0.0668$

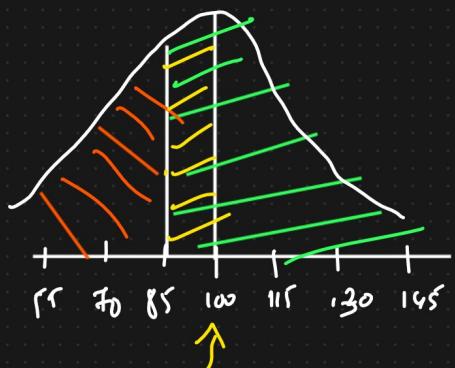
$= 6.6\%$

$$\frac{0.5 - 0.0668}{0.5}$$

Prob In India the average IQ is 100, with a standard deviation of 15. What is the percentage of the population would you expect to have an IQ lower than 85?

Ans)

$$\mu = 100 \quad \sigma = 15$$



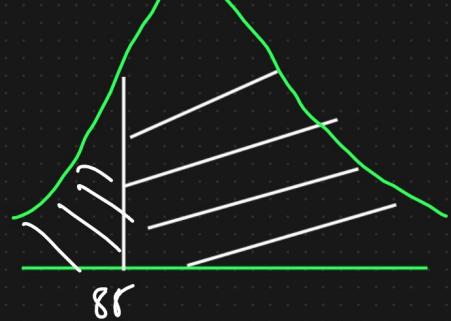
$$\begin{aligned} \textcircled{1} Z\text{-Score} &= \frac{x_i - \mu}{\sigma} \\ &= \frac{85 - 100}{15} \\ &= -1 \end{aligned}$$

Refer Z-table

$$\text{Area under the curve} = \boxed{0.15861}$$

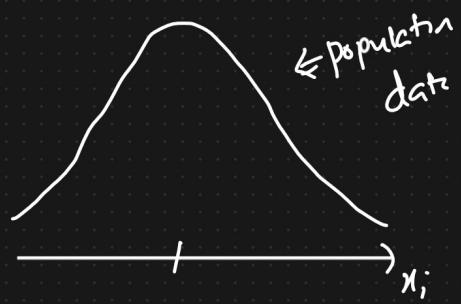
$$\begin{aligned} &0.5 - 0.15861 \\ &= 0.34134. \end{aligned}$$

$$\begin{aligned} \text{Area under the curve } (> 185) &= 1 - 0.15861 \\ &\approx 0.84134. \end{aligned}$$



Central Limit Theorem

$$\textcircled{1} X \sim N(\mu, \sigma) \quad \overbrace{n \text{ samples}}$$



$$\boxed{n=30}$$

Sampling mean

$$S_1 = \{ x_1, x_2, x_3, \dots, x_n \} = \bar{x}_1$$

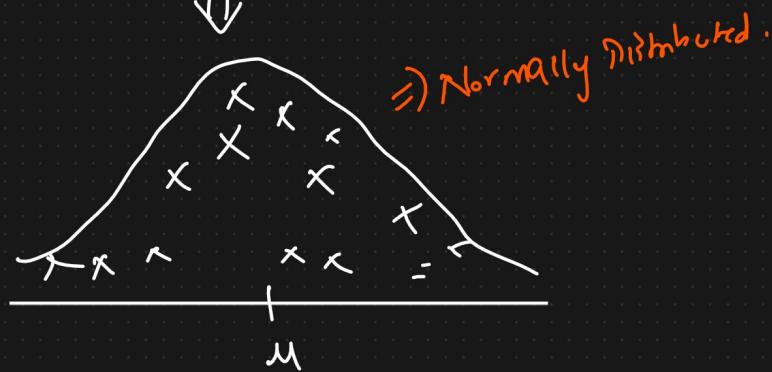
$$S_2 = \{ \dots \} = \bar{x}_2$$

$$S_3 = \{ \dots \} = \bar{x}_3$$

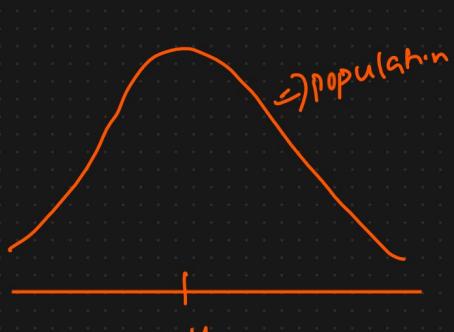
\vdots

$$S_m = \{ \dots \} = \bar{x}_m$$

$\bar{X} = \{\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_m\} \Rightarrow$ Sampling distribution
mean



Important for Interview



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

σ = population std

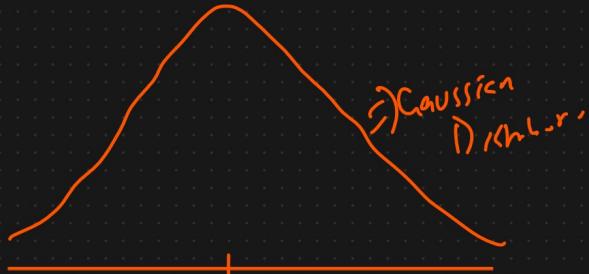
μ = population mean

n = sample size

$$\frac{\sigma}{\sqrt{n}}$$



Sampling distribution of mean



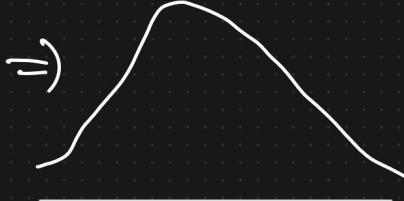
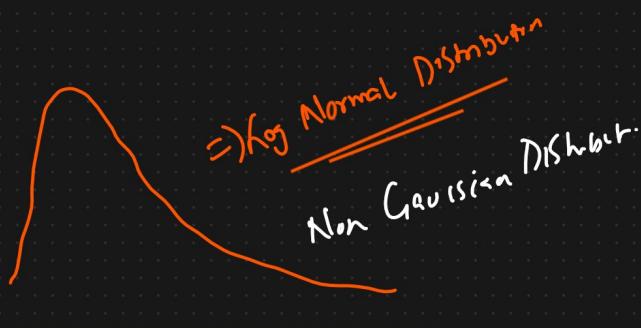
$$\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

n can be any value

Standard Error

$$\boxed{n > 30}$$

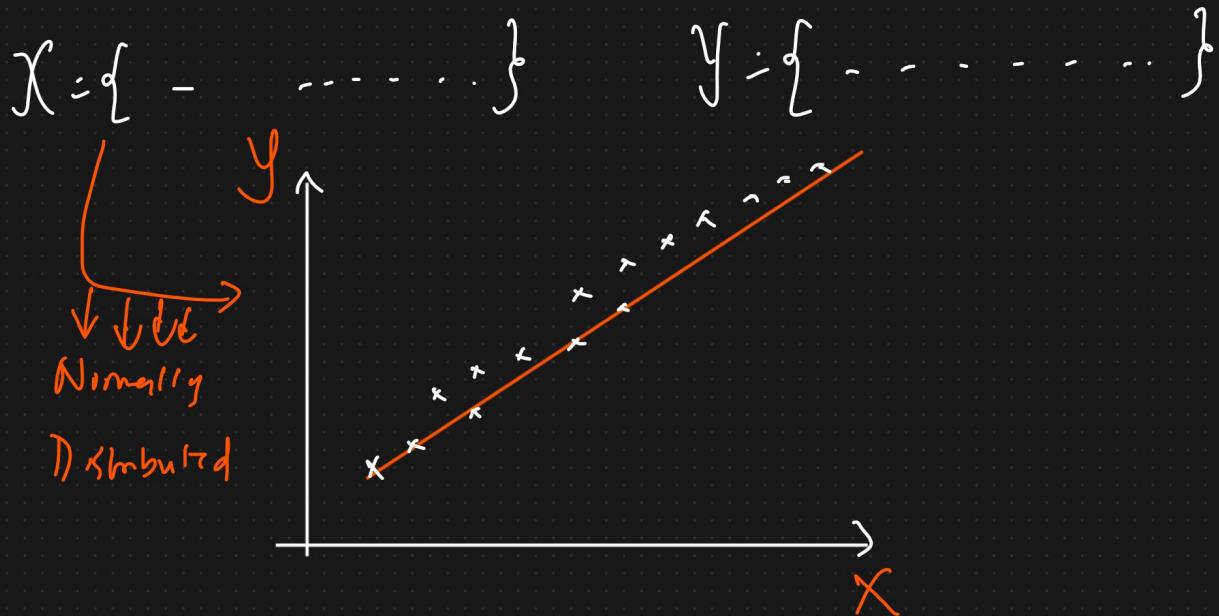
Sampling Distribution mean



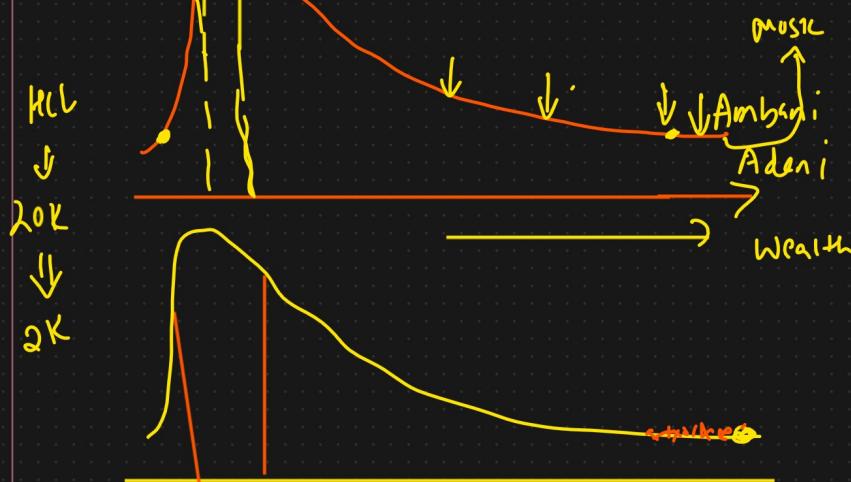
K How to check whether a distribution is Normal or Not

Q-Q plot

Normal Distribution



Log Normal Distribution



Wealth distribution of people in world

