

Statistics

Author: Rohit Mande
Founder and CEO, intrvu.ai

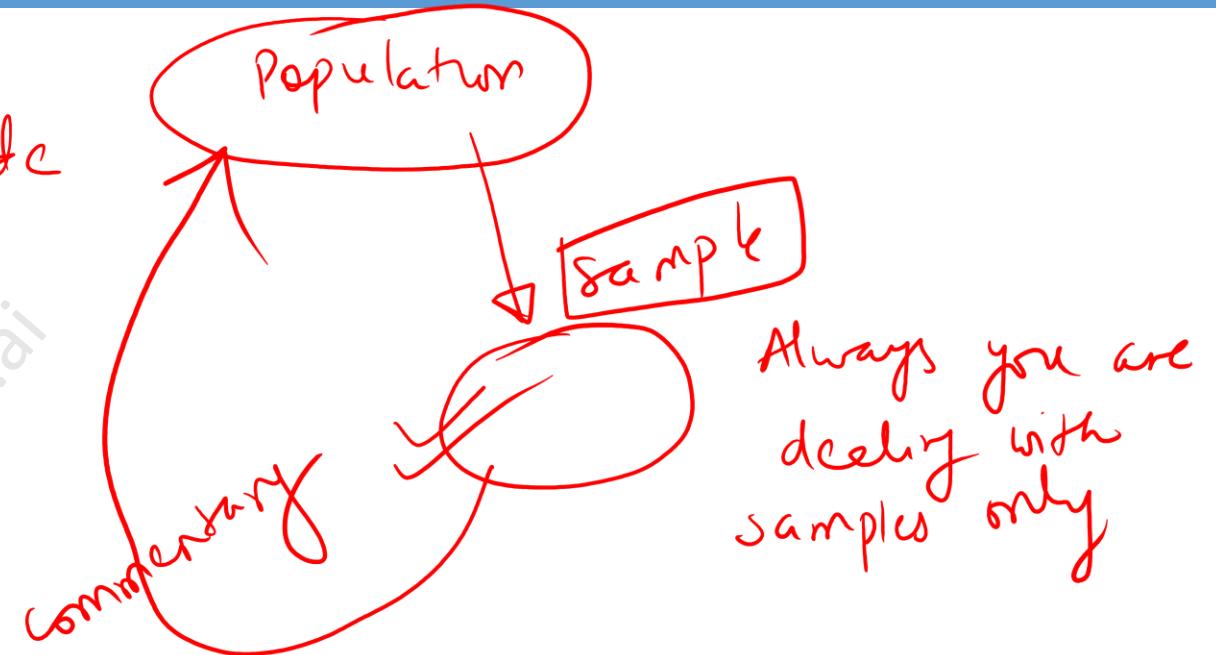


Index

performance of class 10th students in
top schools of India

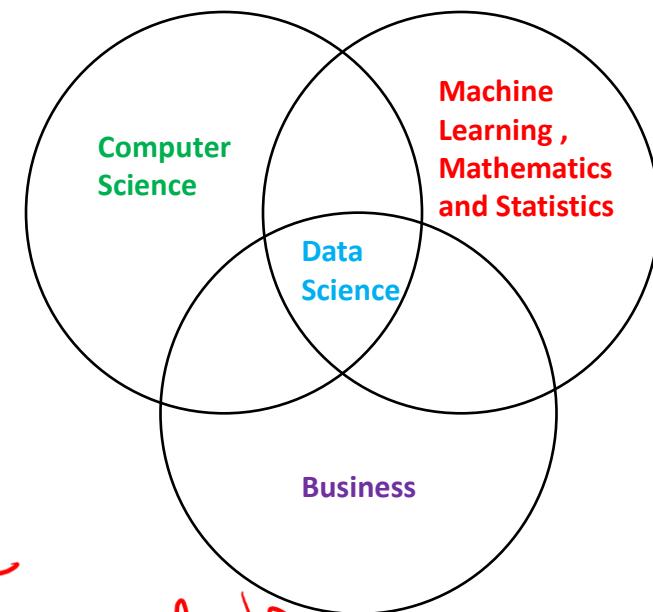
- Data Distribution
- Statistical values mean, median, mode
- Data Normalization
- Missing Value Imputation
- Outlier Detection
- Correlation and Causation

INTTRVU.ai



Why should I learn Statistics?

- Statistics is the discipline that covers analysis, interpretation, and presentation of data
 - Statistics is useful in data science projects for following tasks
 - ✓ Exploratory Data Analysis
 - ✓ Analysing relationship between variables
 - ✓ Missing value imputation
 - ✓ Outlier detection *and treatment*
 - ✓ Hypothesis Testing
- Based on the samples we are commenting on the population*



Q

Approaches for calculating the probability

- 1) MLE (Maximum Likelihood Estimation)
- 2) Bayesian approach

✓ (1) $P(O_1) = \frac{\text{Number of times you observed } O_1}{\text{Total Number of times you performed the EXP}} \rightarrow \infty$

(2) Bayesian

Data Distribution

- ↳ (1) Start with starting knowledge about the phenomenon
- ↳ (2) Experiments
- ↳ (3) Basic knowledge + Experiment → refine your probabilities

Random Variable

$$\begin{array}{r} -10^6, +10^6 \\ \hline -10^6, -10^6+10^6, -10^6+20^6 \\ -10^6+0.5 \times \end{array}$$

Random variable is a variable whose value is unknown e.g., if we measure height of any person on the street then the value, we will get would be any random value (unknown)

- ① **Discrete random variable:** rolling a dice $\{1, 2, 3, 4, 5, 6\}$ $\{1, 2, 3, 4, 5, 6\} \rightarrow \{1, 2, 3, 4, 5, 6\}$
Fix set of values for outcome e.g., tossing a coin

② **Continuous random variable:** 5 ft, 5.1 ft $[1 \text{ ft}, 12 \text{ ft}]$ 5.2 5.25
It can take an infinite (theoretically) number of possible values e.g., salary of a random person selected from Pune

Probability Distribution

→ discrete
→ continuous

Probability distribution is the mathematical function that gives the probabilities of occurrence of different possible outcomes for an experiment

Suppose \underline{X} denotes the outcome of a coin toss ("the experiment"),

As X is discrete random variable (as it takes fix set of values) it is sufficient to specify a function representing probability of each outcome

Probability distribution of X would take the value 0.5 for $X = \text{heads}$,
and 0.5 for $X = \text{tails}$ (for a fair coin).

$$P(X) = \begin{cases} H \rightarrow 50\% \\ T \rightarrow 50\% \end{cases}$$

Experiment

- ① Started tossing the coin for 100 ~~times~~ times } discrete
- ② rolling the dice 100 times
- ③ noting down the salaries of X company

Different possible outcomes for the random variable

$$E_1 = [H, T], E_2 = [1, 2, 3, 4, 5, 6]$$

$$E_3 = [2 \text{ LPA}, 2 \text{ Cr}]$$

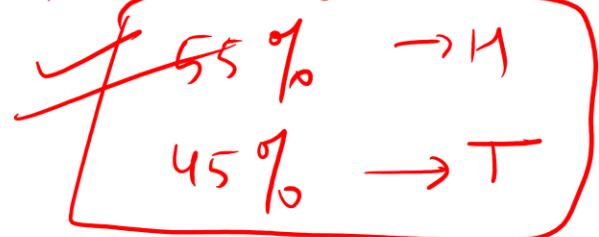
Probability of occurrence

$$E_1: H \rightarrow 55 \quad 55/100 = 0.55$$

$$T \rightarrow 45 \quad 45/100 = 0.45$$

Q1 What is the probability of H occurring? 55%.

Q2 What is the probability distribution for tossing this coin?



$$L \rightarrow \zeta \rightarrow |\zeta|_{\infty}$$

$$2 \rightarrow 16 \rightarrow 1 \cdot 1 \cdot \infty$$

$$3 \rightarrow 17 \rightarrow 17/180$$

$$4 \rightarrow 15$$

5 - 20

$$6 \rightarrow 20 \quad 20/100$$

Bucket size 0.1 LPA

$$2 - 2.1 \text{ LPA} = 3$$

$$21 - 22 \text{ kPa} = 5$$

$$2.2 - 2.3 \text{ kPa} = 7$$

$$10 = 10 \cdot 1 \text{ LPA} = 100$$

200614 - 214 = 1
 INTTRVU

Continuous random variables

1000 employees

Emp1 → 2 APP

Femp 2 → 2.5 LPA

$$\text{Emp}^3 \rightarrow 2.25 \text{ LPA}$$

10 UPA

Emp 10 → 5.1 LIP

\rightarrow 5.5 LPP

$$12 \rightarrow 5.3 \text{ LPR}$$

1

1

$$P(X = x_0) \stackrel{?}{=} \text{Not Defined}$$

\uparrow

exact salary

$$P(x_0 \leq X \leq x_1)$$

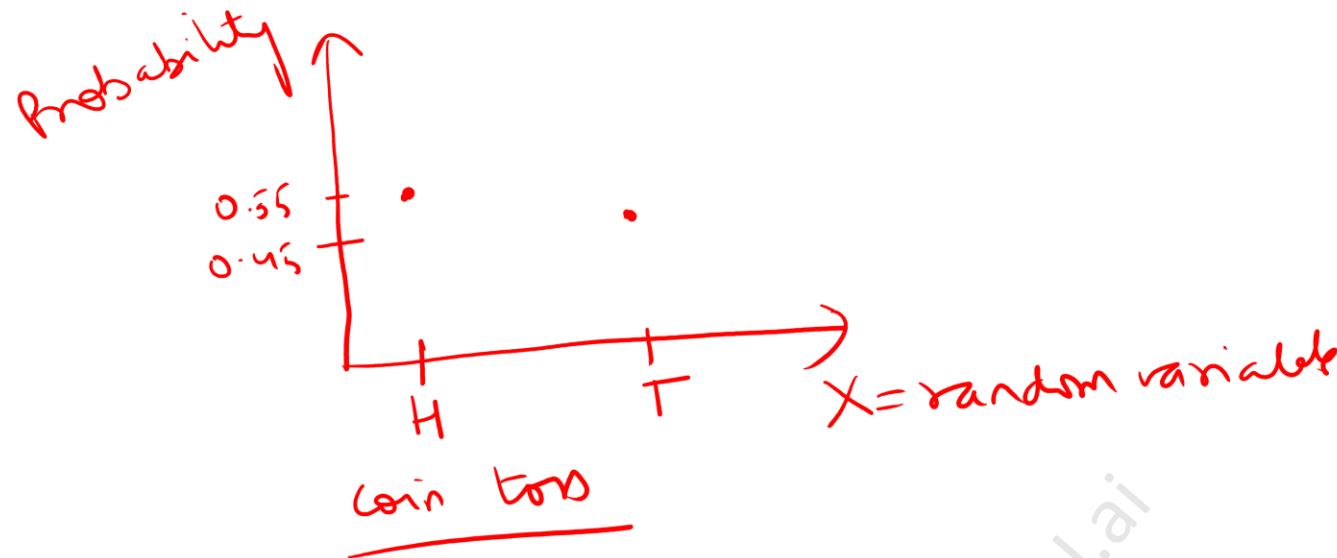
$$P(2 \text{ LPA} \leq X \leq 2.1 \text{ LPA})$$

Probability Distribution

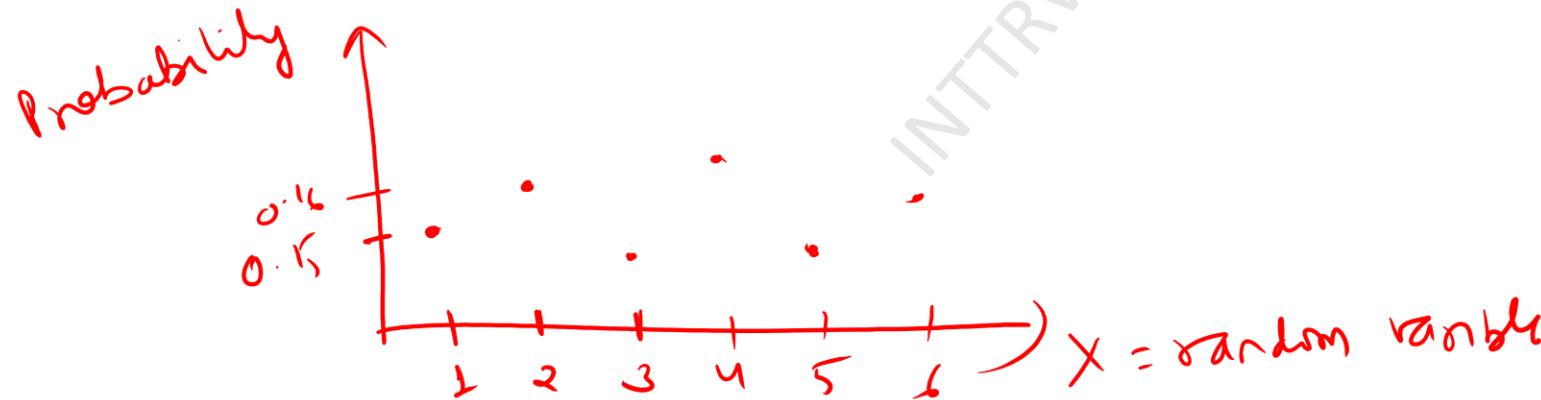
Discrete → Prob. Mass function
Continuous → Probability density function

- Probability distribution is represented using probability density function in case of continuous random variable
- Examples of continuous random variable distributions:
 - Suppose X denotes the outcome of an experiment where we are randomly sampling height of a person from population of city
 - Suppose X denotes the outcome of an experiment where we are randomly sampling average marks obtained by student from all universities

Probability Mass Function for discrete random variables



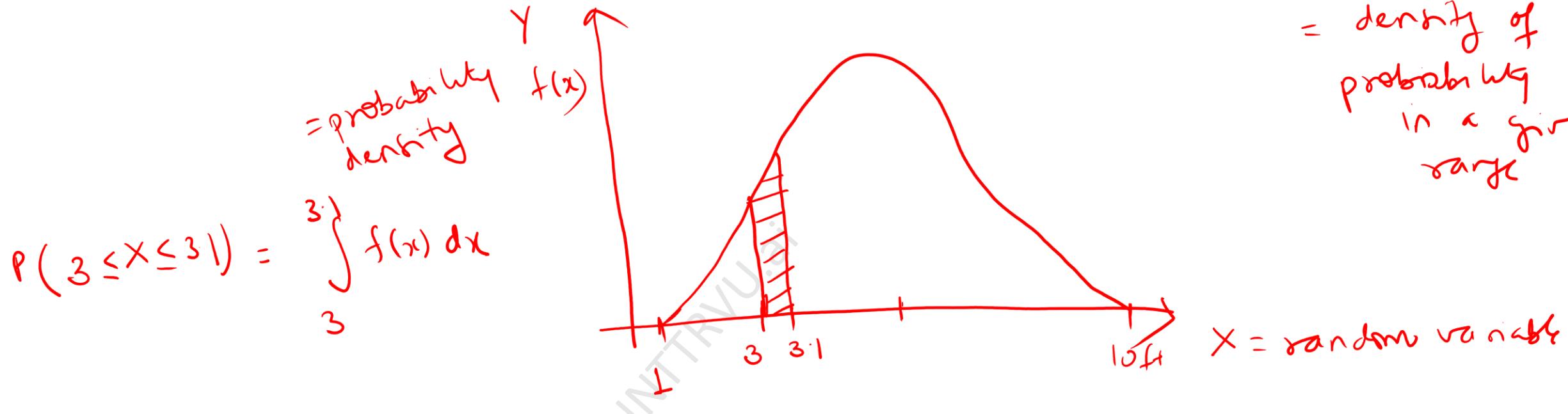
In PMF; we describe the probabilities for each and every possible outcome



For continuous random variable

we can't use PMF

we always talk about pdf (Probability Density function)



γ axis → we don't have the probabilities but we have probability densities

Probability density
= density of probability in a given range

Pdf Probability density function

Observed Height

3 ft
3.2 ft
5 ft
5.6 ft
6 ft
2 ft
2.2 ft

0.1 ft

Buckets	Count
2-2.1	x_0
2.1-2.2	x_1
2.2-2.3	x_2
:	
9.9-10 ft	x_n



reduce the size of each and every bucket 0.01, 0.001 ft
bucket size $\rightarrow 0$

$$\left(\frac{x_0}{\text{Total Count}} \right) \frac{1}{0.1}$$

bucket = 0.1 ft

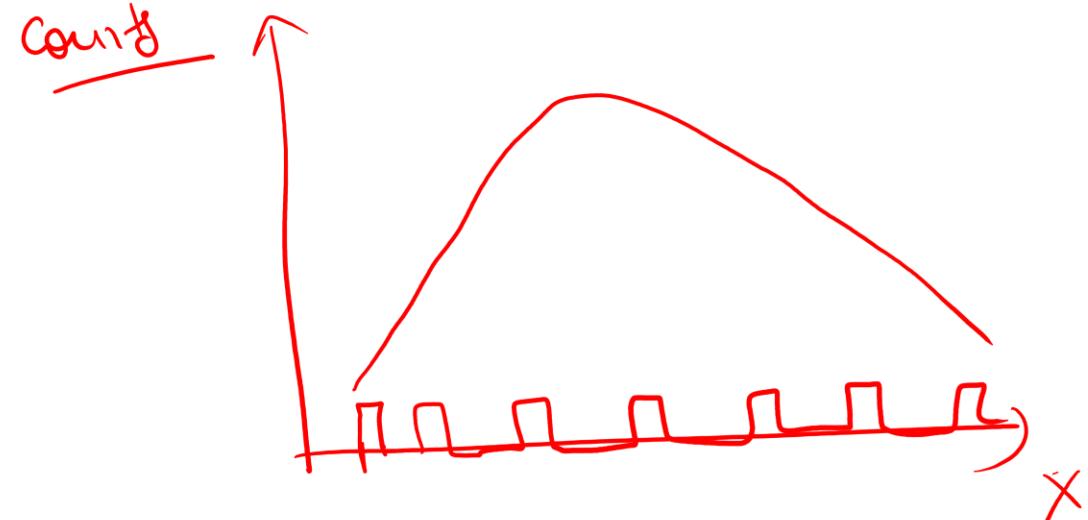
bucket size $\rightarrow 0$

$2 \rightarrow 2.01$

$2 \rightarrow 2.001$

Useless to plot counts or frequencies on Y axis if bucket size $\rightarrow 0$

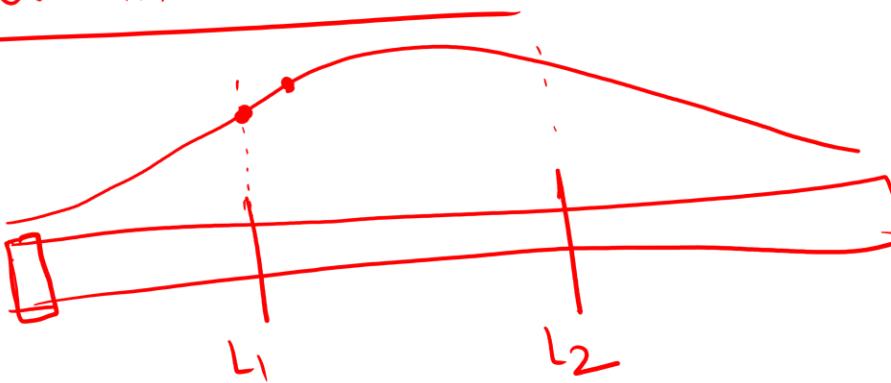
densities $\frac{\Delta x}{\text{density}}$



$$\text{Probability} = \left[\frac{\text{Count}(x_0, x_0 + \Delta x)}{\text{Total count}} \right] \left(\frac{1}{\Delta x} \right)$$

$$\begin{aligned} \text{Probability} &= \int_{y.2LL}^{f(x)} f(x) \Delta x \\ &= \int_y^y (\text{probability density function}) \underline{dx} \end{aligned}$$

Mass of a section in a rod



kg/m

density = $\frac{\Delta M}{\Delta L}$

$0 \rightarrow \Delta L \rightarrow M_1$

$\Delta L \rightarrow 2\Delta L \rightarrow M_2$

$2\Delta L \rightarrow 3\Delta L \rightarrow M_3$

,

$\Delta L \rightarrow \underline{\underline{0}}$

$$M = \int_{L_1}^{L_2} (\text{density}) dL$$

$$\int_{L_1}^{L_2} \frac{\Delta M}{\Delta L} dL$$

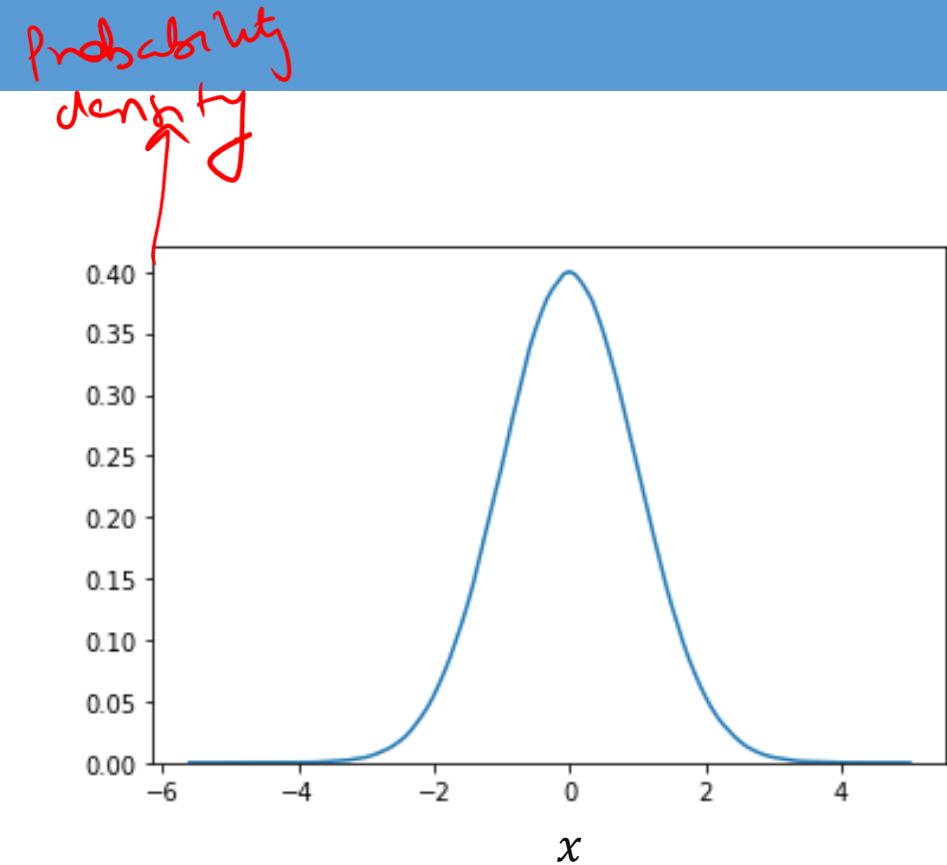
Normal Distribution , bell curve

In statistics, a normal distribution or Gaussian distribution is used to represent probability distribution of a continuous random variable

As sample size increases the shape of sampling distribution becomes more like a normal distribution

Real word continuous random variables are expected to have normal distribution

- Suppose X denotes the outcome of an experiment where we are randomly sampling average marks obtained by student from all universities



$\sigma \rightarrow$ spread of data around μ

Normal Distribution

$\sigma, \mu \rightarrow$ mean
standard deviation

spread of the data
 $\sigma \rightarrow$ 'flatness' of curve

random variable

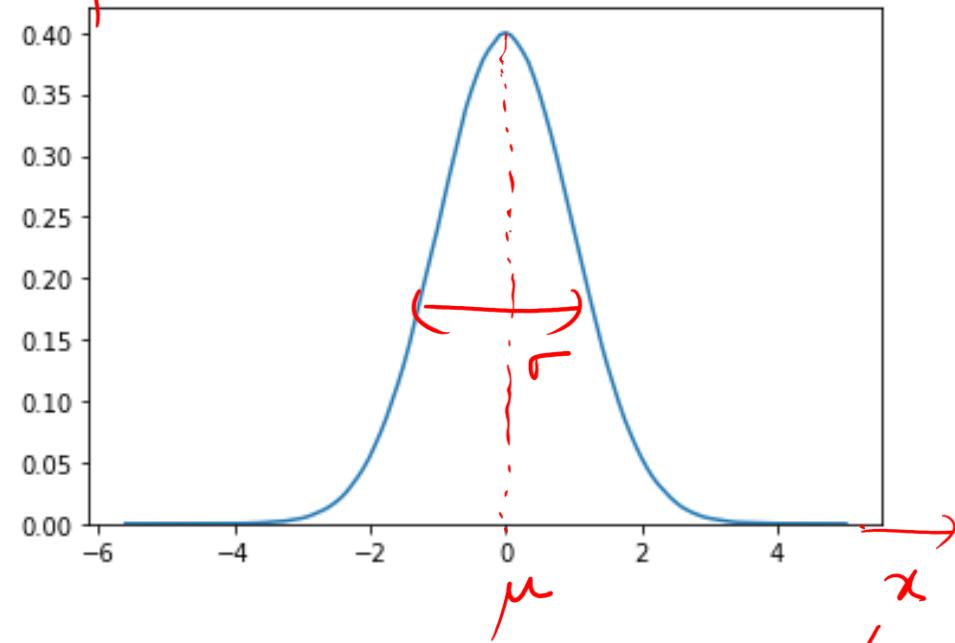
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$f(x)$ - probability density function of random variable x

μ is the mean of the distribution

σ is the standard deviation

probability density



random variable

Normal Distribution

curve is symmetric around μ

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

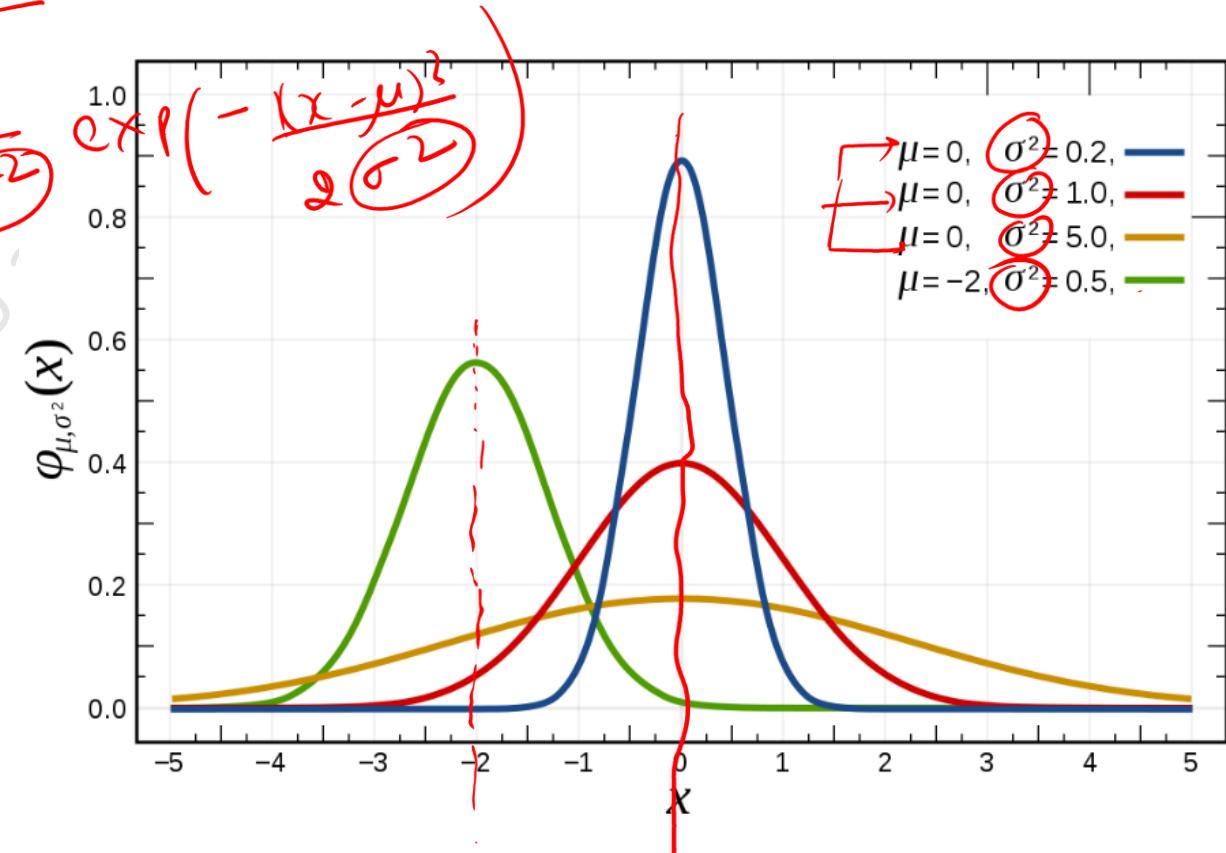
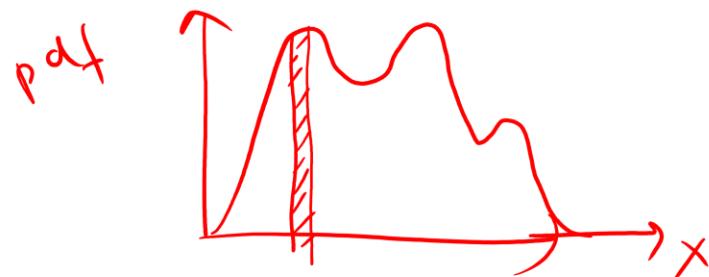
$\approx \frac{1}{\sqrt{2\pi}\sigma}$

$\sigma \uparrow$, fatness \uparrow

$$\boxed{\mu=0, \sigma^2=1}$$

standard normal distribution

Normal distribution with $\mu = 0$ and $\sigma^2 = 1$ is called as standard normal distribution



Height = $[h_0, h_1, h_2, \dots, h_{999}]$

-1, 1] \rightarrow 0

① $\mu = \frac{h_0 + h_1 + h_2 + \dots + h_{999}}{1000}$

descriptive statistic $-5, 5] \rightarrow 0$

② spread of data points around μ

(ft) $h_0 \rightarrow h_0 - \mu = \text{diff}_0 (\text{ft})$
(ft) $h_1 \rightarrow h_1 - \mu = \text{diff}_1 (\text{ft})$
(ft) $h_2 \rightarrow h_2 - \mu = \text{diff}_2 (\text{ft})$
...
(ft)² $h_{999} \rightarrow h_{999} - \mu = \text{diff}_{999} (\text{ft})$

Variance = $\frac{(\text{diff}_0)^2 + (\text{diff}_1)^2 + \dots + (\text{diff}_{999})^2}{1000}$

standard deviation = $\sqrt{\text{Variance}}$

$\frac{(\text{ft})}{\sqrt{1000}}$

Central Limit Theorem

Population

→ sample 1 → μ_1

→ sample 2 → μ_2

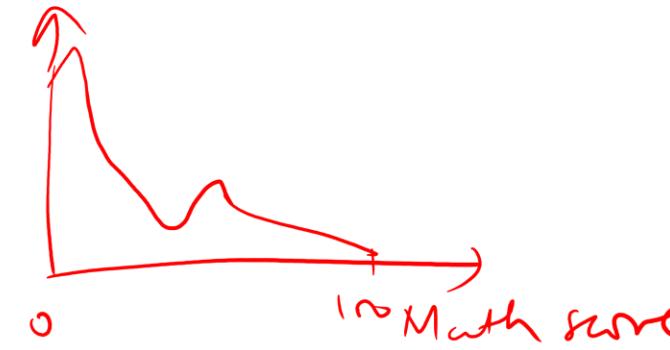
:

→ sample 100 → μ_{100}

→ sample 101

:

→ sample 250



→ pdf → Normal distribution as we keep
on increasing sample number

Central limit theorem

The central limit theorem states that if you take sufficiently large number of (>30) samples from a population, the samples mean will be normally distributed, even if the population isn't normally distributed.

E.g., Sample salary of 1000 people living in Pune and plot its distribution

$[h_0, h_1, h_2, \dots, h_{999}] \rightarrow$ 1000 height values

$$\mu = \frac{h_0 + h_1 + h_2 + \dots + h_{999}}{1000}, \text{ median}$$

Step 1 sort the data

$[h'_0, h'_1, h'_2, \dots, h'_{999}]$

Descriptive Statistics

Mode

most occurring data point ; data point with highest probability

$[1, 1, 2, 3, 3, 3, 4, 5]$

$1 \rightarrow 2$

$2 \rightarrow 1$

$3 \rightarrow 3$

$4 \rightarrow 1$

$5 \rightarrow 1$

Statistical Values

Step 2 To identify the middle data point

Step 1 $[3, 1, 2] \downarrow$ median
 $(1, 2, 3)$

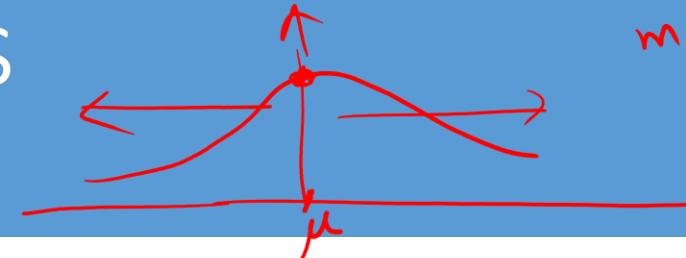
$[4, 1, 3, 2]$

$(1, [2, 3], 4)$

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

Mean, median, mode, variance, standard deviation

Statistical values



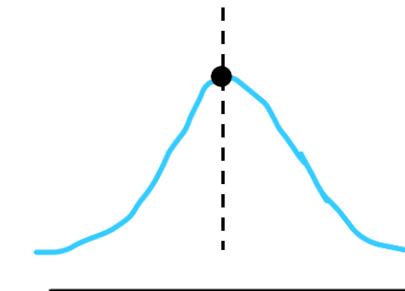
mean = median = mode

Mean (μ): It is the arithmetic average value of the dataset

central tendency

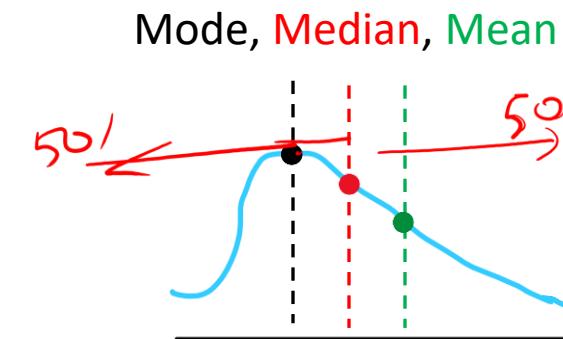
Median: Median is the middle value in sorted dataset. It is calculated by sorting the data in ascending order and then calculating the value at middle index

Mean, Median, Mode

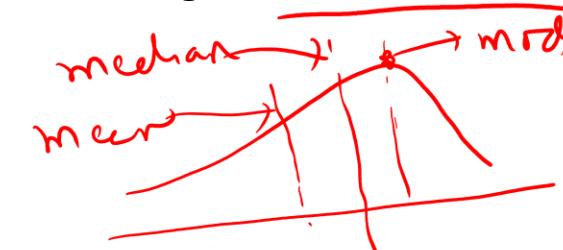


Mode: It is the most occurring value in the dataset

symmetric distribution



Right skewed distribution



Statistical values

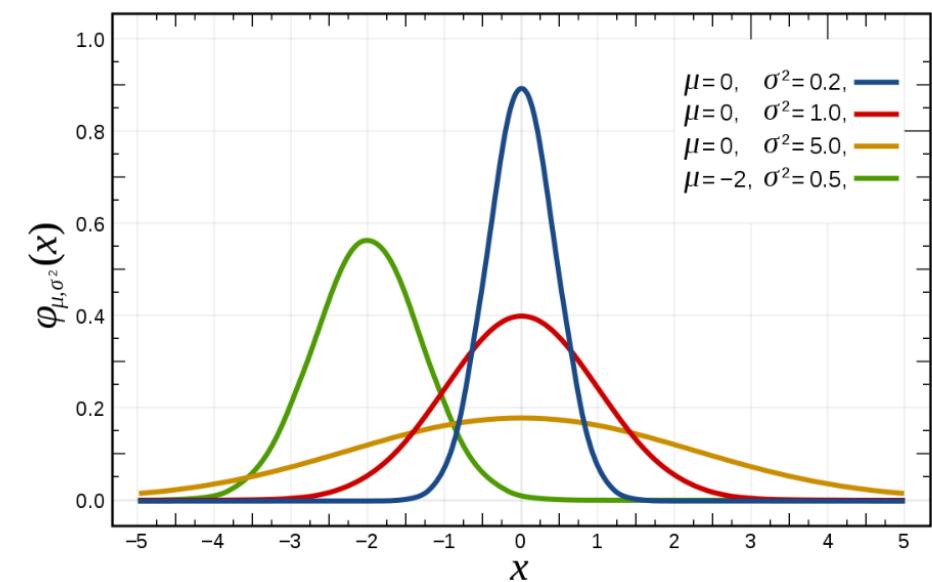
Deviation: It is the difference between the observation and mean of the sample population

$$\text{Deviation} = \underline{x_i} - \mu$$

Variance: It is mean of squares of deviation

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

Standard Deviation (σ) : It is the square root of the variance



Height	deviation
h_0	$h_0 - \mu = d_0$
h_1	$h_1 - \mu = d_1$
h_2	$h_2 - \mu = d_2$
.	.
h_{999}	$h_{999} - \mu = d_{999}$

$$\begin{array}{|c|c|} \hline & \text{d} \\ \hline d_0 & d_0^2 \\ d_1 & d_1^2 \\ d_2 & d_2^2 \\ \vdots & \vdots \\ d_{999} & d_{999}^2 \\ \hline \end{array}$$

Variance: $\frac{d_0^2 + d_1^2 + d_2^2 + \dots + d_{999}^2}{1000}$ $\mu = \frac{\sum h_i}{1000}$

↓
spread in the
data

$$\begin{array}{l} [-1, 0, 1] \rightarrow 0 \\ [-5, 0, 5] \rightarrow 0 \end{array}$$

spread of the data

I have 1000 deviation value, how do I summarize them?

$$\frac{d_0 + d_1 + d_2 + \dots + d_{999}}{1000}$$

$$\begin{bmatrix} -1 - 0 = d_0 = -1 \\ 0 - 0 = d_1 = 0 \\ 1 - 0 = d_2 = 1 \end{bmatrix}$$

$$\frac{d_0 + d_1 + d_2}{3} = 0$$

Variance

(Dataset 1)

$$\text{Data} = [-2, -1, 0, 1, 2]$$

Step 1 $\mu = \frac{(-2) + (-1) + 0 + 1 + 2}{5} = 0$

Step 2

Data	Deviations
-2	$d_0 = (-2) - \mu = -2$
-1	$d_1 = (-1) - \mu = -1$
0	$d_2 = 0 - \mu = 0$
1	$d_3 = 1 - \mu = 1$
2	$d_4 = 2 - \mu = 2$

$$\begin{aligned}
 \text{Variance} &= \frac{d_0^2 + d_1^2 + d_2^2 + d_3^2 + d_4^2}{5} \\
 &= \frac{(-2)^2 + (-1)^2 + 0^2 + 1^2 + 2^2}{5} \\
 &= \frac{4 + 1 + 0 + 1 + 4}{5} = \frac{10}{5} = 2
 \end{aligned}$$

$$\text{standard deviation} = \sqrt{\text{variance}}$$

$$= \underline{\underline{\sqrt{2}}}$$

Variance

Data set 2

$$\text{Data} = [-10, -5, 0, 5, 10]$$

$$\mu = 0 = \frac{-10 + (-5) + 0 + 5 + 10}{5}$$

Data	deviation
-10	-10 - 0 = -10
-5	-5
0	0
5	5
10	10

$$\text{Variance} : \frac{1}{5} \left[(-10)^2 + (-5)^2 + 0^2 + 5^2 + 10^2 \right]$$

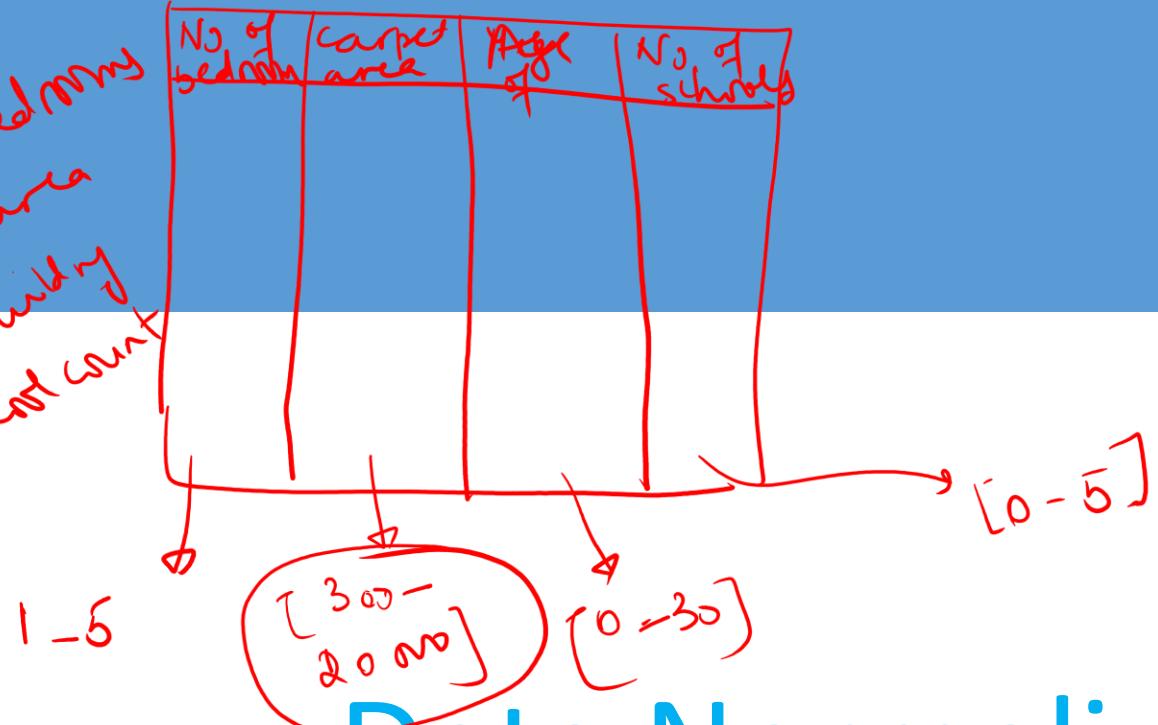
$$= \frac{1}{5} \left[100 + 25 + 0 + 25 + 100 \right]$$

$$= \frac{250}{5} = 50$$

$$\text{Std} = \sqrt{50}$$

Dataset

μ_1 for #bedrooms
 μ_2 for carpet area
 μ_3 for #schools
 μ_4 for #schools



house price prediction data.

range of the values of different features becomes a hindrance to the model

Data Normalization

No. of bedrooms

$$[0, 1]$$

carpet area

$$[0, 1]$$

δ_x

$$[0, 1]$$

No. of schools [0, 1)

$$\text{DP1} \rightarrow (2, 1000, 5, 2)$$

$$\text{DP2} = (4, 2000, 10, 4)$$

$$\sqrt{2^2 + (1000^2) + 5^2 + 2^2}$$

Data Normalization

- Normalization is used to scale the data column (feature) such that all features fall in the same range e.g., 0.0 to 1.0 or to transform the data into distribution with zero mean and standard deviation/variance 1
- It is a recommended step for some machine learning algorithms to get correct results e.g., k-means clustering

$$DF_1 \rightarrow [0.5, 0.6, 0.2, 0.3]$$
$$DF_2 = [0.52, 0.8, 0.3, 0.5]$$

$$\sqrt{(0.2)^2 + 0.2^2 + 0.1^2 + 0.2^2}$$

Data Normalization

→ Z Score
→ min - max

Z-score normalization:

This method transforms data such that mean will be zero and standard deviation will be 1 for transformed data

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

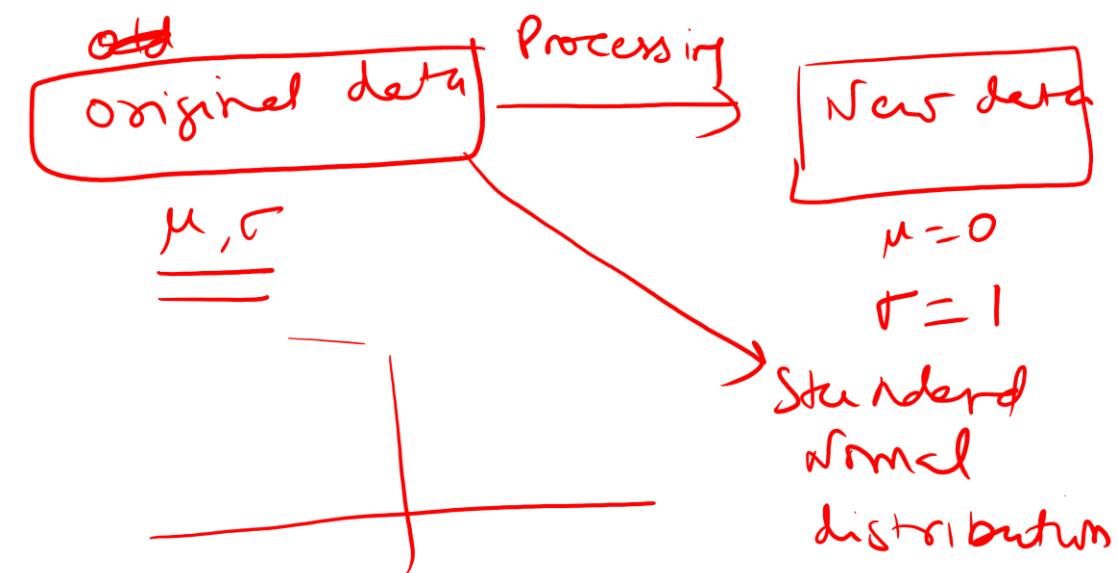
z is the new value of feature after applying normalization

x is original value of the feature

μ is mean value

σ is the standard deviation

INTTRVU.ai



Step 1

μ, σ for original data

Step 2

$$z_i = \frac{x_i - \mu}{\sigma}$$

x_i = i^{th} data point

height

Original height	Normalized height
h_0	$\frac{h_0 - \mu}{\sigma}$
h_1	$\frac{h_1 - \mu}{\sigma}$
h_2	$\frac{h_2 - \mu}{\sigma}$
\vdots	
h_M	$\frac{h_M - \mu}{\sigma}$

$\underbrace{\hspace{100px}}$ Normalised column

μ, σ for height

Goal was that new column should have $\text{new } \mu = 0$ and $\text{new } \sigma = 1$

$$\text{new } \mu = \frac{h_0 - \mu}{\sigma} + \frac{h_1 - \mu}{\sigma} + \frac{h_2 - \mu}{\sigma} + \dots + \frac{h_M - \mu}{\sigma}$$
$$= \frac{(h_0 - \mu) + (h_1 - \mu) + (h_2 - \mu) + \dots + (h_M - \mu)}{\sigma (M+1)}$$

$$= \frac{(h_0 + h_1 + h_2 + \dots + h_M) - (M+1)\mu}{\sigma (M+1)}$$

$$= \frac{1}{\sigma (M+1)} \left[\frac{(h_0 + h_1 + h_2 + \dots + h_M)}{M+1} - \mu \right]$$

Data Normalization

z score normalization
 $\hookrightarrow \mu=0, \sigma=1$

Min-Max normalization:

$$[3.00, 10.00] \rightarrow [0, 1]$$

min *max*

This method transforms data to the range of 0 to 1 (minimum value be 0 and maximum value will be 1 after transformation)

$$z_i = \frac{x_i - \min(x)}{\max(x) - \min(x)}$$

z_i is the new value of feature after applying normalization

x_i is original value of the feature

$\min(x)$ is the min value of feature

$\max(x)$ is the max value of feature

height	Normalized height
h_0	$\frac{h_0 - h_i}{h_j - h_i} > 0$
h_1	$\frac{h_1 - h_i}{h_j - h_i} > 0$
h_2	$\frac{h_2 - h_i}{h_j - h_i} > 0$
h_i	$\frac{h_i - h_i}{h_j - h_i} = 0$
h_j	$\frac{h_j - h_i}{h_j - h_i} = 1$
h_m	$\frac{h_m - h_i}{h_j - h_i} > 0$

Step 1 find min height
and max height

$h_i \rightarrow$ minimum height

$h_j \rightarrow$ maximum height

$[0, 1]$
 new
minimum new
maximum

Data = $[-5, -1, 0, 1, 5, 6, 10]$

min = -5, max = 10

$$z_i = \frac{x_i - \text{min}}{\text{max} - \text{min}}$$

-5	$\frac{-5 - (-5)}{10 - (-5)} = 0$
-1	$\frac{-1 - (-5)}{10 - (-5)} = \frac{4}{15}$
0	$\frac{0 - (-5)}{10 - (-5)} = \frac{5}{15}$
1	$\frac{1 - (-5)}{10 - (-5)} = \frac{6}{15}$
5	$\frac{5 - (-5)}{10 - (-5)} = \frac{10}{15}$
6	$\frac{6 - (-5)}{10 - (-5)} = \frac{11}{15}$
10	$\frac{10 - (-5)}{10 - (-5)} = \frac{15}{15} = 1$

Objective : real world may have missing data points

ML models can handle complete dataset only

put "some value" in the place of missing value

	Age	bill	hrs	steep	Exercy	Y
I1		v				1
I2		-				0
I3		x				1
I4		x				0

"\n" =

Missing Value Imputation

"Informed value" → we will decide what to put in missing value with the help of existing data

Missing value Imputation

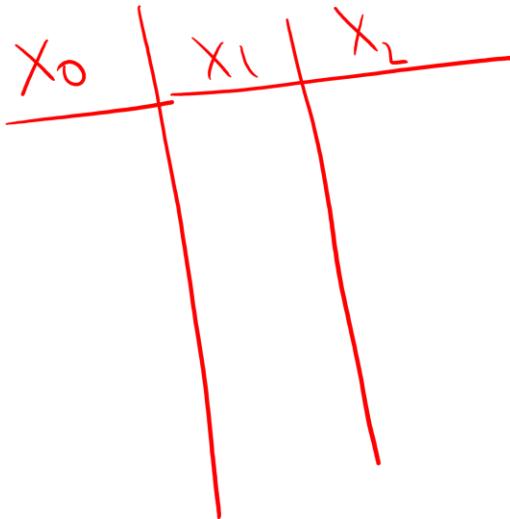
- In real world dataset we might get features with missing values
- To use those features in ML models we have to impute the missing values
- **Missing value imputation for numerical feature**
 - Mean value of feature – mean is sensitive to outliers
 - Median value of feature
- **Missing value imputation for categorical feature**
 - Mode value of feature



Temperature
23
25
21
25
26
29
24
27
28

Missing value
Impute by mean: 25.3

Mean method of missing value imputation



$\frac{x_0}{10}$
 20
 30
 15 → x
 40
 45

$$\bar{x} = \frac{10+20+30+15+40+45}{6}$$

Median method

$$\bar{x} = \text{median}(10, 15, \boxed{20, 30}, 40, 45) = \frac{20+30}{2} = \underline{\underline{25}}$$

Mean vs Median (x_0)

① Mean is affected by outliers

1, 2, 3, 4, 5, 10000

$$\mu = \frac{1+2+3+4+5+10000}{6}$$

Numerical

$\frac{10+1}{2}$

9.89

5.5

$P(X = x_0) = \text{Not Defined}$

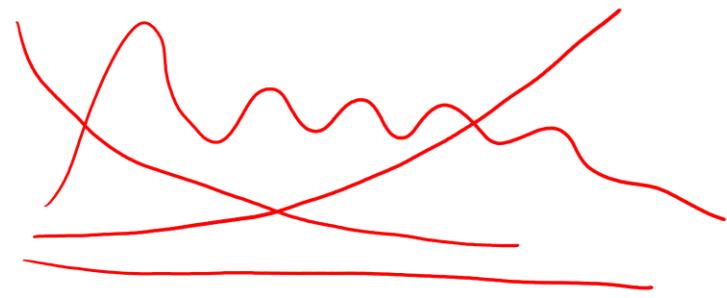
$P(x_1 \leq X \leq x_2) = \underline{\underline{=}}$

median = (1, 2, 3, 4, 5, 10000)

$$\frac{3+4}{2} = \underline{\underline{3.5}}$$

② Data distribution

plot the histogram for x_0



Missing value imputation for categorical features

X_1	0	1	2	3	4	0	1	1	1	1
'S'										
'M'										
'L'										
'X1'										
'XXL'										
'S'										
'M'										
'M'										
'M'										
'M'										

\Rightarrow Mode is 'M'

\Rightarrow we will input the missing value with mode

$$\text{Mean} = \frac{0+1+2+3+4+0+1+1+1+1}{10} = 1.4$$

$$0+1 = 1$$

$$'S' + 'M' \neq 'M'$$

0, 1, 2, 3, 4
 'S', 'M', 'L', 'X1', 'XXL'
 'M', 'F'
 0 + 2

\rightarrow categorical

$$'L' - 'M' = \underline{\underline{M}}$$

$| > 0$ -2
 $(F) > (M) \rightarrow$ makes no sense
 0.5 0.6

Missing value imputation

x ₀	x ₁	x ₂	x ₃	x ₄
✓	✓	✓	✓	✓
✓	✓	✓	✓	✓
✓	✓	✓	✓	✓
✓	✓	✓	✓	✓

100,000 total number of rows

100 rows have missing data

It is going to change the data distribution

80%
data
is missing

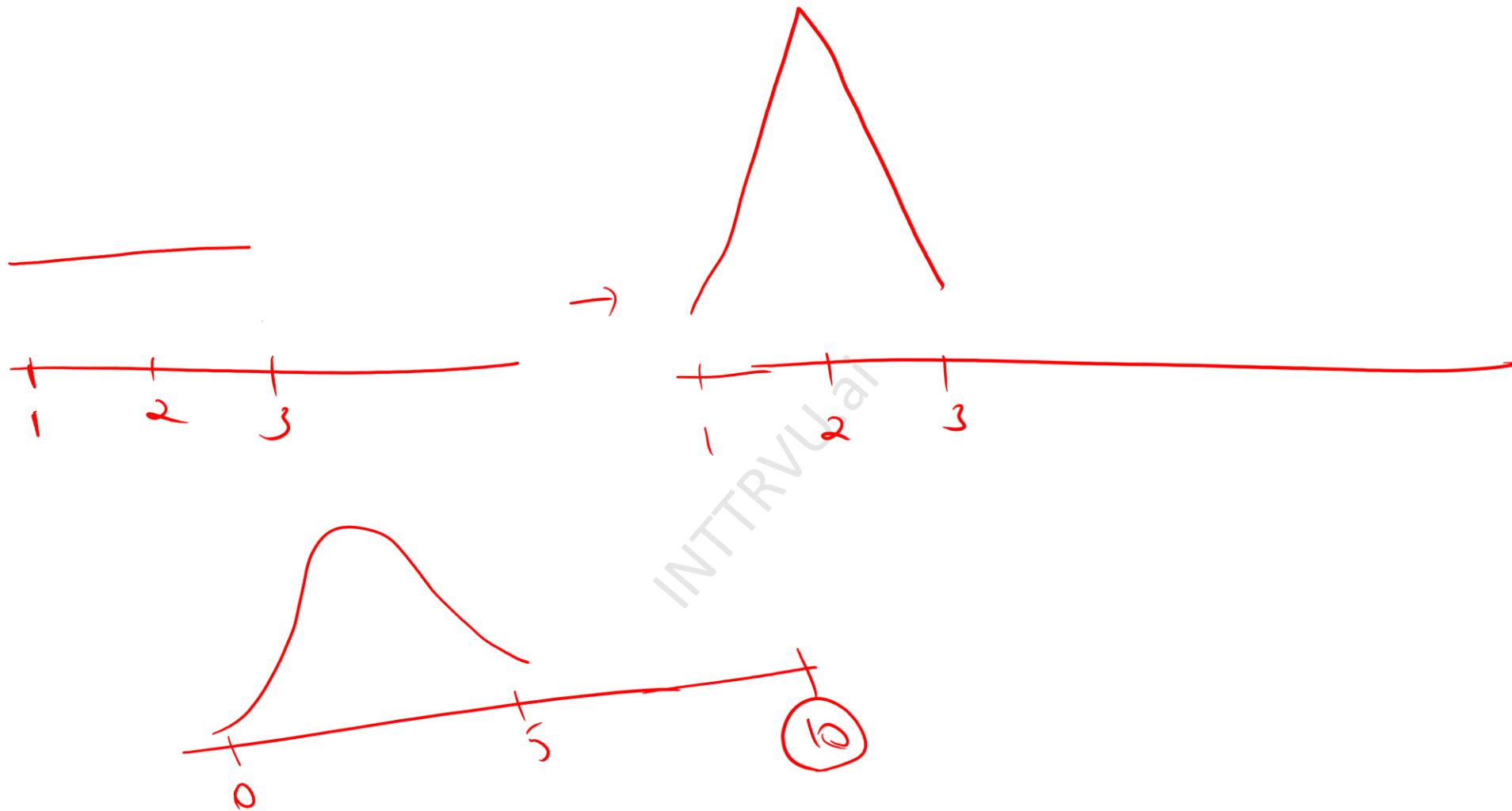
x ₀	x ₁	x ₂	x ₃	x ₄
✓	✓	✓	✓	✓
✓	✓	✓	✓	✓
✓	✓	✓	✓	✓
✓	✓	✓	✓	✓

No missing value imputation

→ we simply drop the column

data distribution

[1, 2, 3, 2, 2, 2, 1, 2, 2, 2, 2, 2]



Outlier \rightarrow value which is different from the remaining values in dataset

$[1, 2, 3, 4, 5, 1000]$

$[1, 2, 3, 4, 5, 10]$

Why do we care about outlier? it changes the underlying data distribution (meters)

Outlier Detection

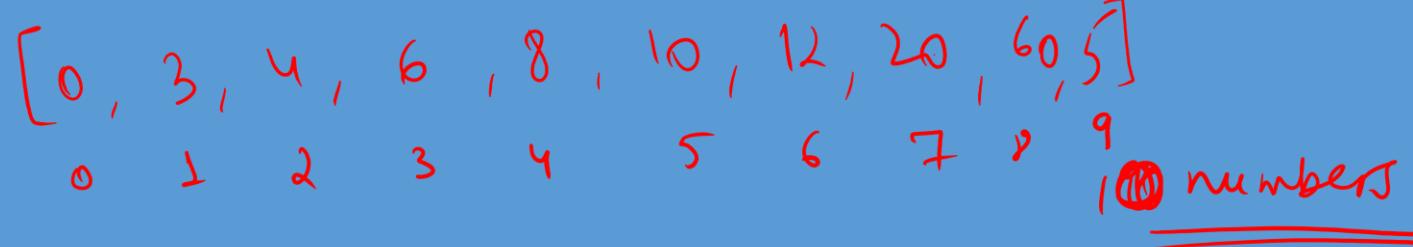
"centimeters \rightarrow error"

min	max
1.3	2.4
<u>240</u>	<u>—</u>

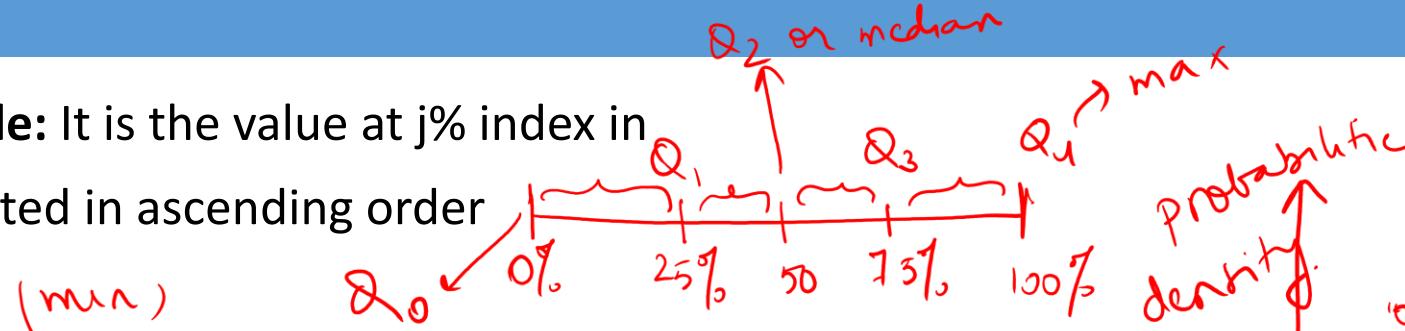
Student1 $\frac{400}{500} \times 100 = 80\%$

Student2 $\frac{900}{1000} \times 100 = 90\%$

Percentiles



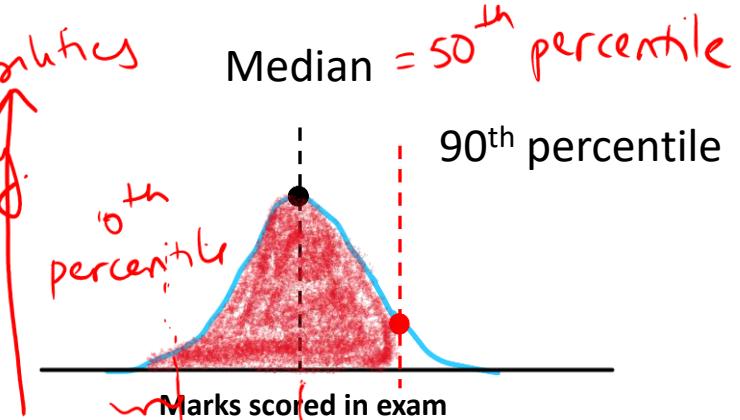
Percentile: It is the value at $j\%$ index in data sorted in ascending order



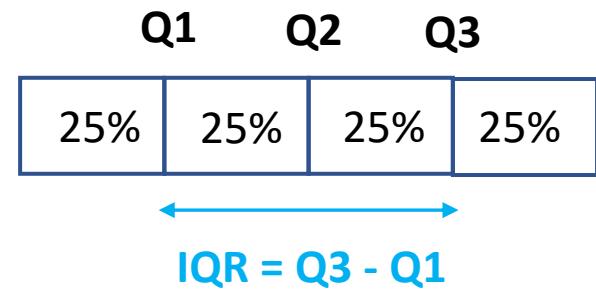
Lower Quartile (Q1): It is the value at 25% index in data sorted in ascending order

Upper Quartile (Q3): It is the value at 75% index in data sorted in ascending order

Inter Quartile Range: $IQR = Q3 - Q1$



Gathering the marks in an ascending order
10,000



Array 1 $[3, 8, 10, 20, 30, 40, 45, 50, 55, 70]$ electrical
0 1 2 3 4 5 6 7 8 9
engineering

$$\begin{array}{c} \uparrow \\ 50\% \\ \hline \end{array}$$

$$\frac{5}{10} = \underline{\underline{50\%}}$$

$$\frac{10}{10} = \underline{\underline{100\%}}$$

Array 2 $[8, 10, 30, 35, 85]$ mechanical engineering
0 1 2 3 4


Percentiles are useful when we try to compare the ranks

Percentiles are percentages
for array indices

Percentiles are "percentages of array indices"

Step 1 Take the array and sort it

Input array $\rightarrow [0, 8, \cancel{2}, \cancel{4}, \cancel{6}, 10, 20, \cancel{1}, 30, 15]$ 50th percentile
 $10 \times 0.5 = 5$

$[0, 1, \cancel{2}, \cancel{4}, 5, 6, 8, 10, 15, 20, 30]$ → indices

25th percentile

$10 \times 0.25 = \underline{\underline{2.5}}$ interpolation method
2nd number → 1st index
or
3rd number 2nd index

↳ answer will be number present at 2.5th index

Outlier Detection

$$Q_1, Q_3$$

$$IQR = Q_3 - Q_1$$

$$\textcircled{1} \text{ threshold } L = \text{Lower threshold} = Q_1 - 1.5^* IQR$$

$$\textcircled{2} \text{ Upper threshold} = Q_3 + 1.5^* IQR$$

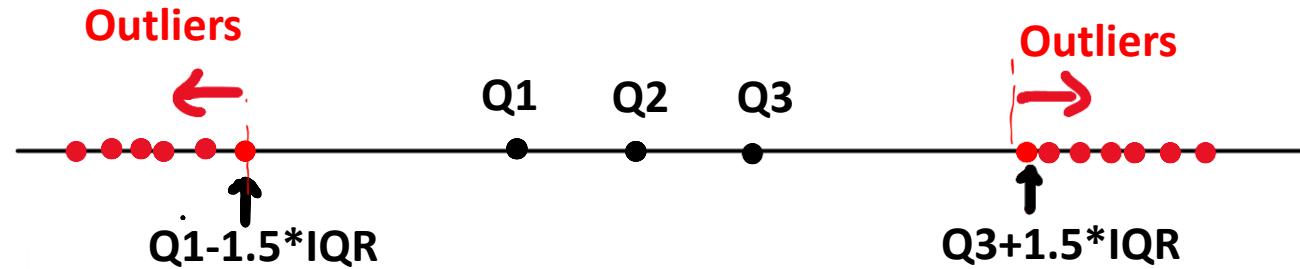
Outlier is an observation which lies at an abnormal value from other values in sample population

IQR is useful in identifying the outliers

good values
↑

$$[Q_1 - 1.5^* IQR, Q_3 + 1.5^* IQR] \text{ ---}$$

Potential outliers lie outside the range: [Q1-1.5*IQR , Q3+1.5*IQR]





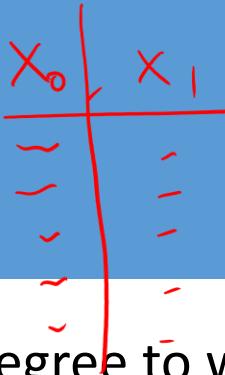
Potential outliers

Bad outliers

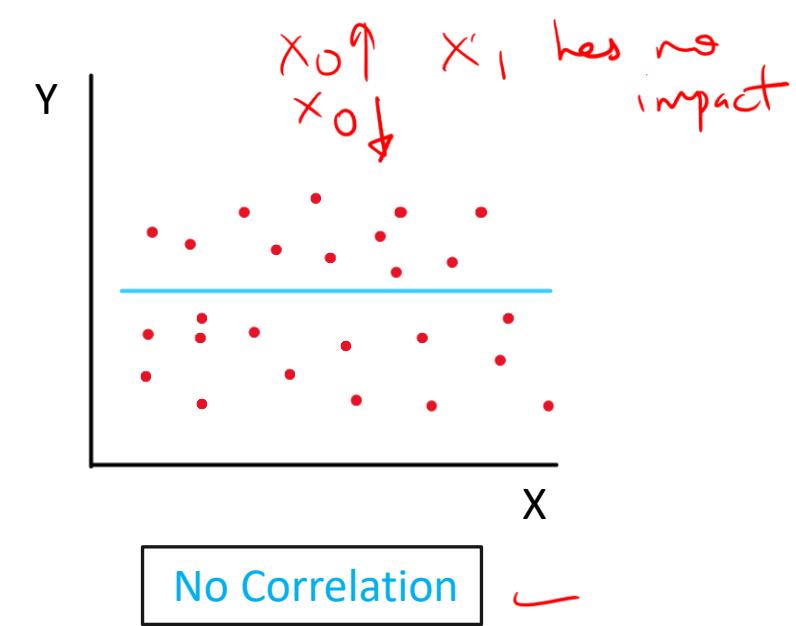
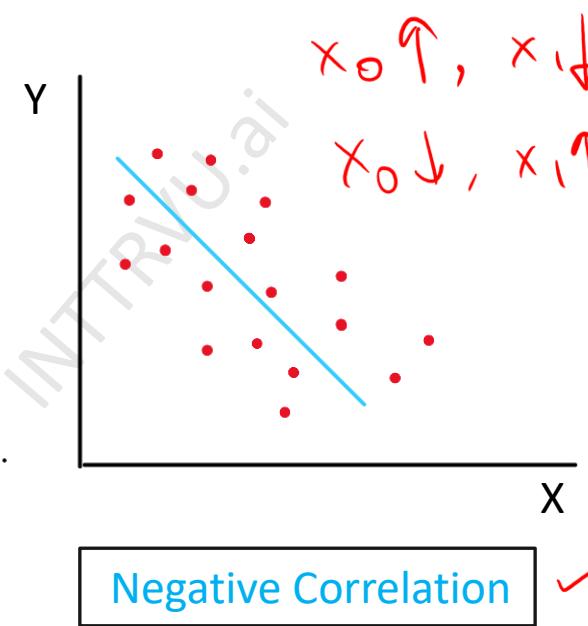
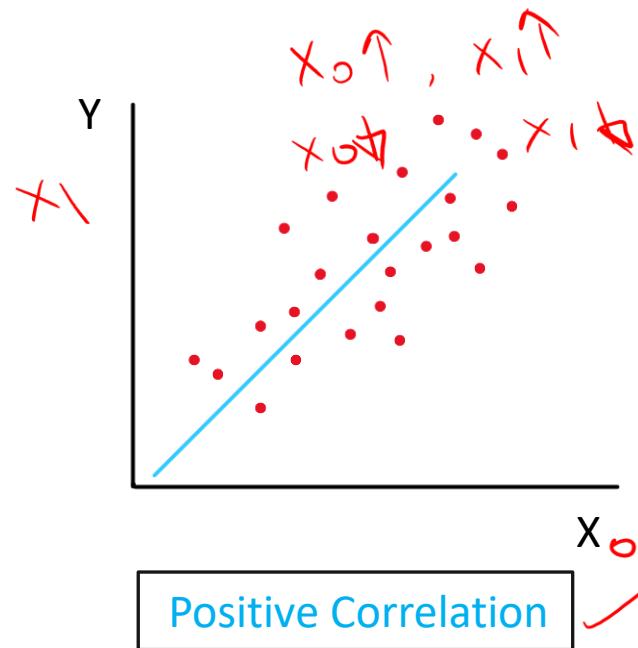
Good outliers

Correlation and Causation

Correlation



Correlation indicates the degree to which a pair of variables are linearly related



① Pearson's Correlation Coefficient [-1 to 1]

It is calculated by dividing the covariance of the two variables by the product of their standard deviations.

It is in the range of -1 (perfect negative linear relationship) to 1 (perfect positive linear relationship)

Values near 0 indicate that there is no correlation between the variables

$$\check{p_{X,Y}} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

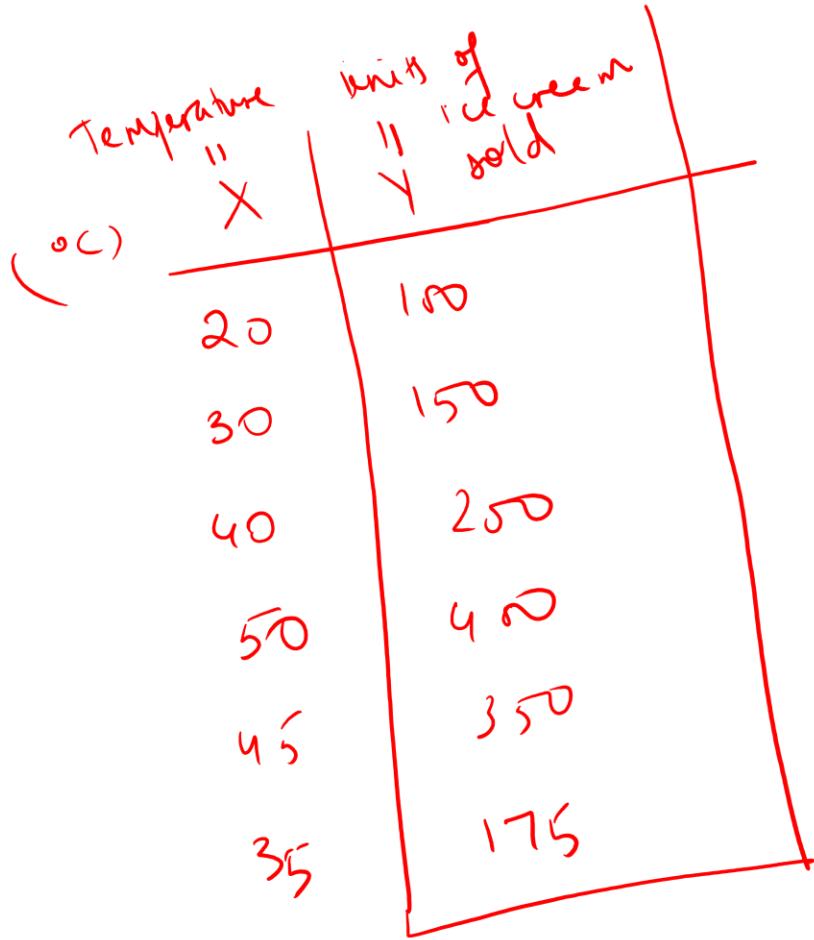
-1 to 1

Covariance is a measure of the joint variability of two random variables (e.g., Higher values of X generally correspond to higher values of Y)

High covariance will give high correlation coefficient as well

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

σ_X, σ_Y = standard deviations of X and Y columns



Expectation / mean values in Y column

$$E[(X - \mu_X)(Y - \mu_Y)]$$

values in X column

mean of X column

mean of Y column

co-varying

co-variance
= variation together

$$E(X) = \frac{20 + 30 + 40 + 50 + 45 + 35}{6}$$

Temp X	Units of ice cream sold Y	$X - \mu_X$	$Y - \mu_Y$	$(X - \mu_X)(Y - \mu_Y)$
20	100	-16.67	-129.16	16.67×129.16
30	150	-6.67	-79.16	6.67×79.16
40	200	3.33	-29.16	-3.33×29.16
50	400	13.33	170.84	13.33×170.84
45	350	8.33	120.84	8.33×120.84
35	125	-1.67	-54.16	1.67×54.16

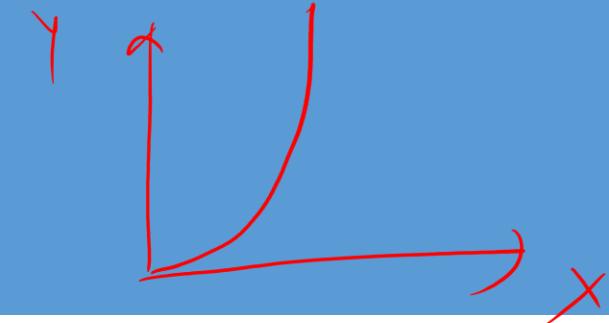
$$\mu_X = 36.67$$

$$\mu_Y = 229.16$$

$$E[(X - \mu_X)(Y - \mu_Y)] = \text{Covariance}$$

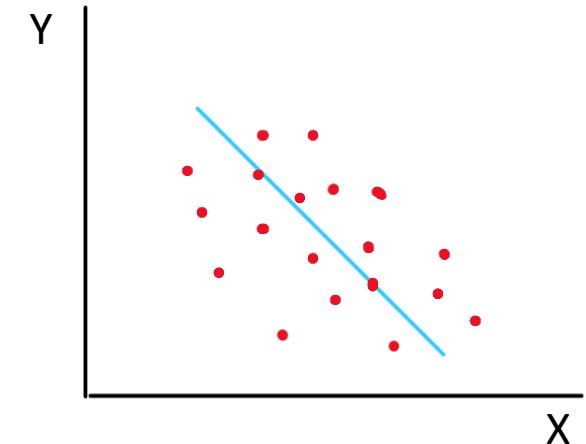
Correlation

x_0	x_1	x_2	x_3	y



Correlation for feature selection:

- Correlation is useful in finding numerical features with predictive power which can be used in regression model
- Low correlation does not mean that feature is not useful as correlation covers only linear relationship
- While building machine learning model it is recommended to use the feature in model and evaluate its importance if the feature is useful in that business domain



Negative Correlation

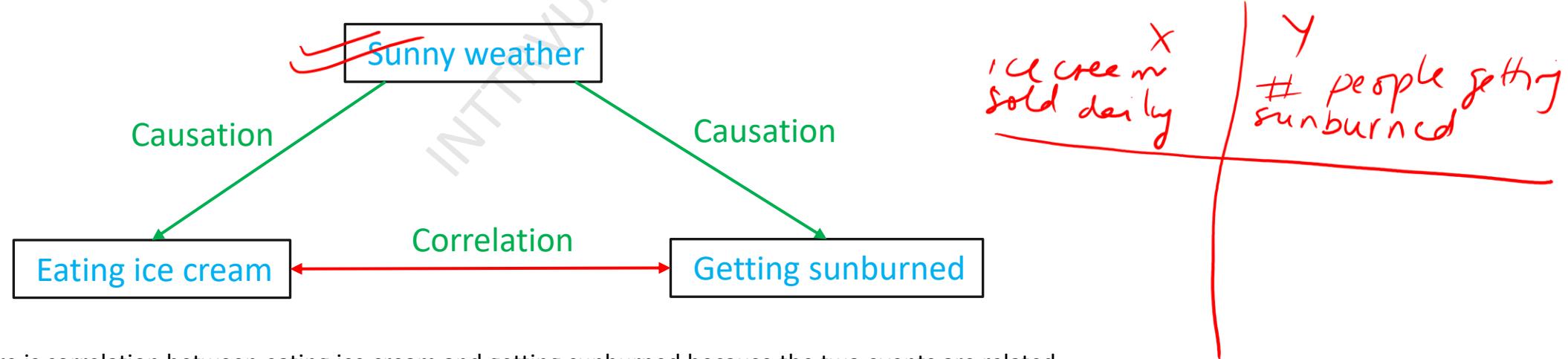
X and Y = $P = 0.99$ but does X cause Y ? causation
 Y causes X .

Causation

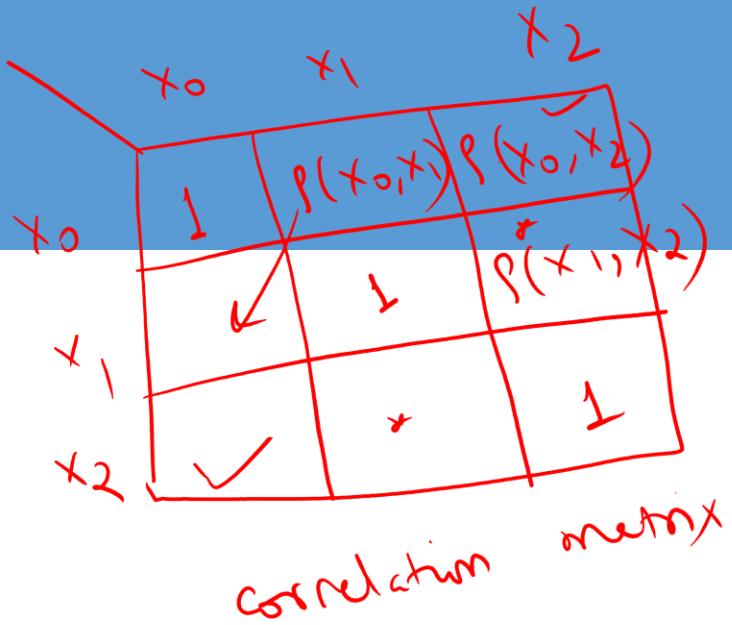
While correlation and causation can exist at same time, correlation does not imply causation

Causation means one thing causes another — e.g., action A causes outcome B.

On the other hand, correlation is simply a relationship where action A relates to action B—but one event doesn't necessarily cause the other event to happen.



There is correlation between eating ice cream and getting sunburned because the two events are related. But neither event causes the other. Instead, both events are caused by something else—sunny weather

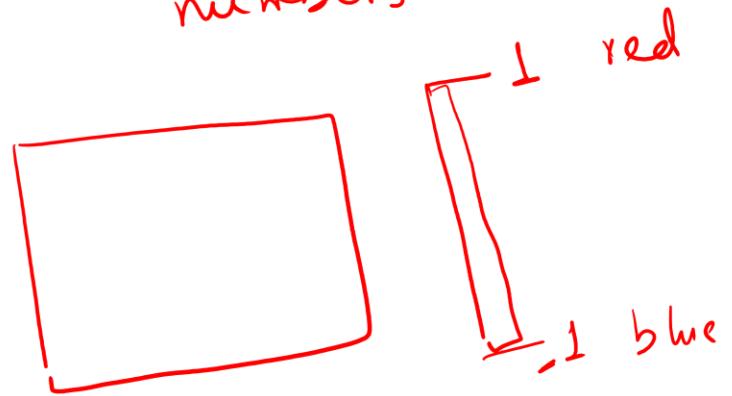


heat map

x_0 = weight
 x_1 = BP
 x_2 = cholesterol

1 feature from row
 1 feature from column

color mapping with the numbers



Questions?

Using ML mode for missing value imputation

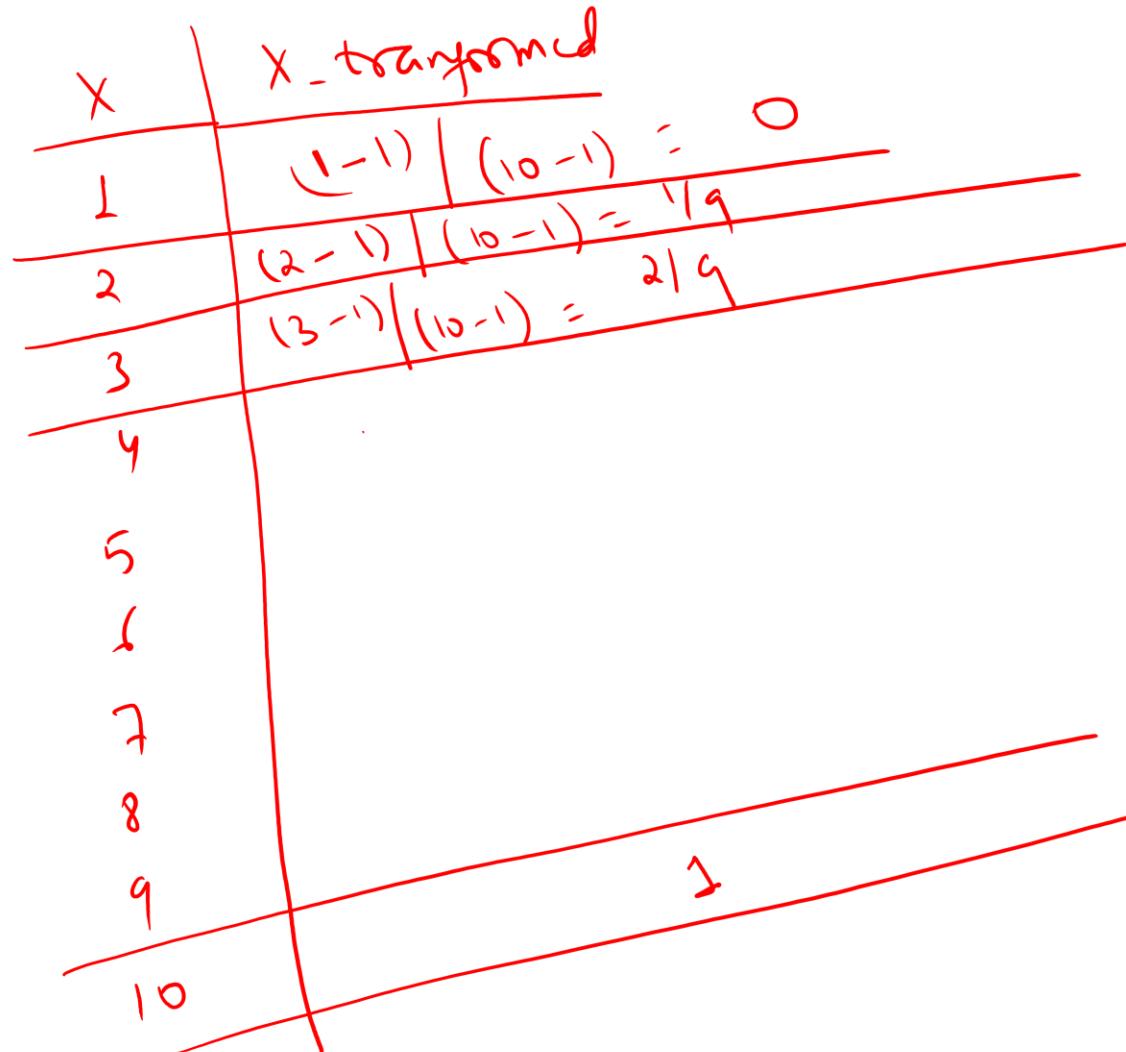
Data 1 having no missing values \rightarrow train a model x_0, x_1, x_2, x_3
(400 rows) \hookrightarrow features and x_4 as the predict column

→ Data 2 with missing values of Xy
(100 rows)

A hand-drawn diagram illustrating a machine learning process. On the left, the input data $(1, 2, 3, y)$ is shown. An arrow points from this data to a red-bordered box labeled "Trained ML model". Another arrow points from the box to the right, leading to the label "predicted X_1 value".

Normalization

sklearn → Min Max Scalar



x_i

$$z_i = \frac{x_i - \min}{\max - \min}$$

min = 1
max = 10

fit
"method"

$$z_i = \frac{x_i - \min}{\max - \min} \rightarrow \text{transform}$$

"method"

$$\text{fit_transform} = \text{fit} + \text{transform}$$

Standard Scales

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

μ, σ → fit

$\frac{x - \mu}{\sigma}$ → transform

