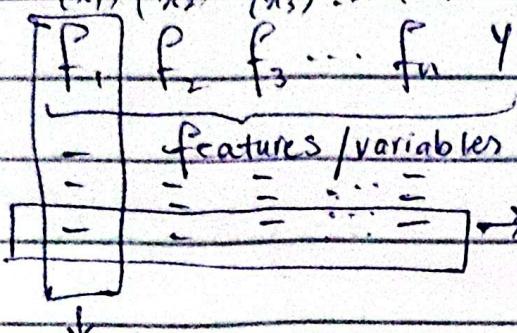


Exploratory Data Analysis

↳ Univariate

one variable

$(x_1), (x_2), (x_3), \dots, (x_n)$  → label/target



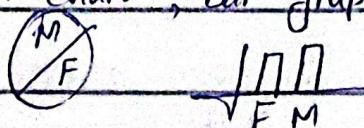
height → higher analysis

feature

Univariate → pick one feature for analysis

Descriptive analysis → we provide summary

↳ Visual analysis → Pie chart, bar graph



→ Python code

```
import matplotlib.pyplot as plt
```

```
fig1 = plt.figure() row location } Row wise fill
```

```
fig1.add_subplot(Plot 2, 2, 1) }
```

```
plt.scatter('x-val', column 'y-val')
```

- Scatterplot ( $X, Y$ )

grid type:

which = "major" → draw <sup>on</sup> major points

- PIE PLOT (Categorical Data > Variable)
  - plt.pie(value, label="label-name") (By default = %, can also do count)
- Erich (Beautify) - Term in DS to add values
  - output = '%.2f' %' → 2 decimal place + % sign
  - auto percentage
  - starting and ending format
- startangle = 90 → default = 0°
  - value to start the plotting
- explode = [0, 0, 0, 0, 0, 0] → default = 0 means no explode
  - wedgeprops Number of pics/values in data
- BAR PLOT - (Preferred Categorical Data)
  - Labels are not sorted in sort(). It only sort the values count.
    - ① either sort the labels with values count
    - ② or use sort(sort=False), then it will not sort and show the default values.
- By Default = count, also do count with slightly change

BAR Chart → categorical Data (Freq. Distribution)

VS

• HISTROGRAM (Used for ~~discrete~~ discrete value)

~ Use bin for discrete data to

$$\text{bin} = \text{int}((\max(\text{values}) - \min(\text{values})) / 2)$$

- Density plot → For bin (very difficult to evenly distribute of real data)

distribution → (Seaborn is used for density())

PDF → to show frequency distribution

import seaborn as sns

① hist draw: dist. plot

ax = sns.distplot(data=values, kind='kde')

distribution plot

Kernel density estimate

① hist → histogram

② ecdf → empirical (preferred) ( $\frac{\text{Discrete}}{\text{Continuous}}$ ) variables

↳ used to show normal distribution

① kde → (second preference) ( $\frac{\text{Discrete}}{\text{Continuous}}$ )

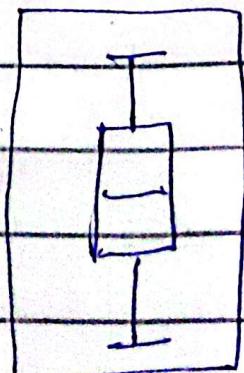
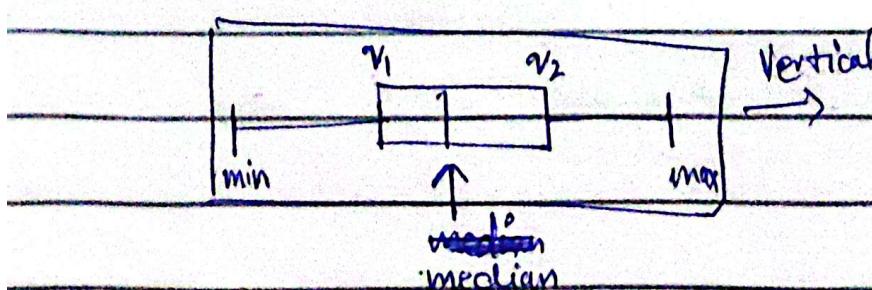
Normal distribution

→ Bell shape → if data is normally dist.

Then shape shows very good shape

• BOX PLOT

aka 5 number summary



- Discrete  
Continuous

Friday

IDS

23-1

## Bivariate Data Analysis

- Scatter, Line, Count, Box, Swarm Plots

① Numerical Data    ② Categorical Data

Correlation → Increase of independent

cause increase / decrease in dependent var.

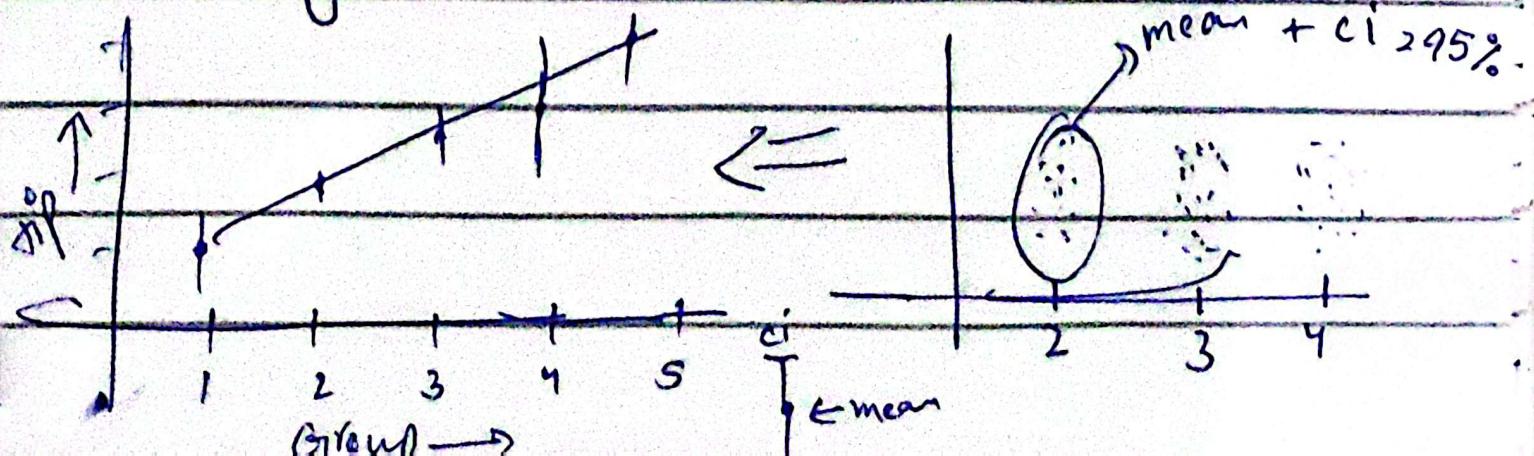
Causality → How much change in dependent  
is caused by independent

sns. relplot() → relation plot

sns. regplot() → add regression line

- for continuous data ( $x$ -estimator) mean  
or average of data at a point is

use along with confidence interval

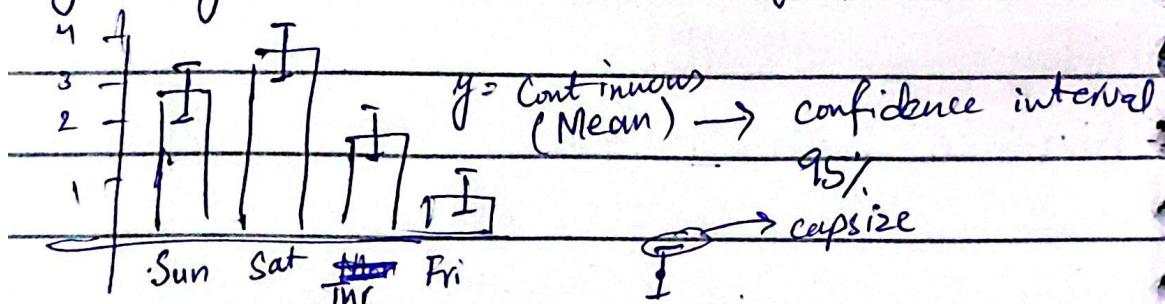


`sns. lmplot()` → Linear plot

Categorical  $\xrightarrow{\text{on}}$  against  $\rightarrow$  Continuous

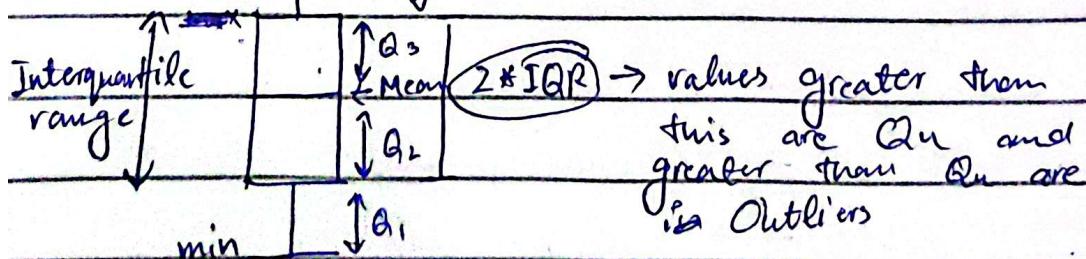
① `Countplot()` → countplot used for categorical with continuous data. e.g. how many male, female paid tips.

E.g. Day and customer ~~fun~~ group

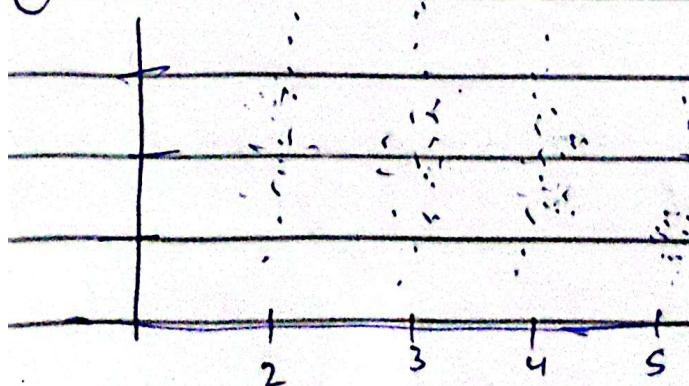


X = Categorical

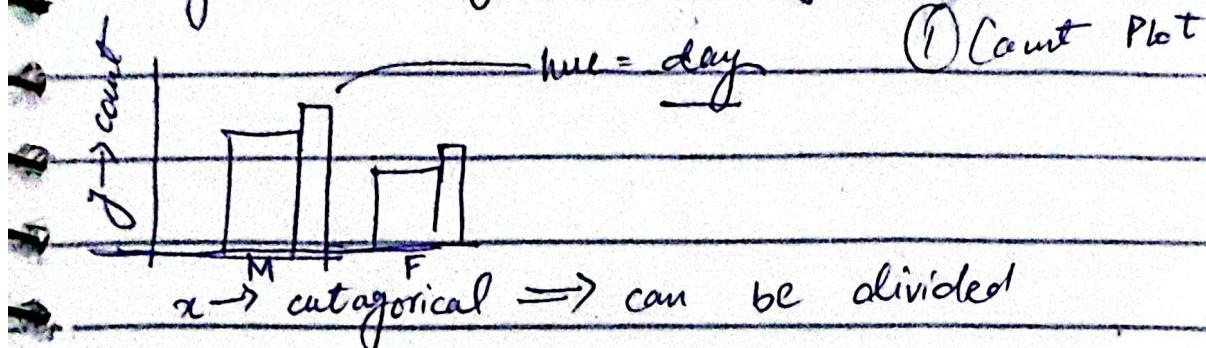
② Boxplot  $\circlearrowleft$  → Outliers  
max  $\leftarrow Q_1$  summarize the data



③ Swarm Plot



Categorical  $\rightarrow$  against  $\rightarrow$  categorical



② Distribution / Density Plot (Continuous vs Categorical)

can be raw

↓ column

bar

No specific condition

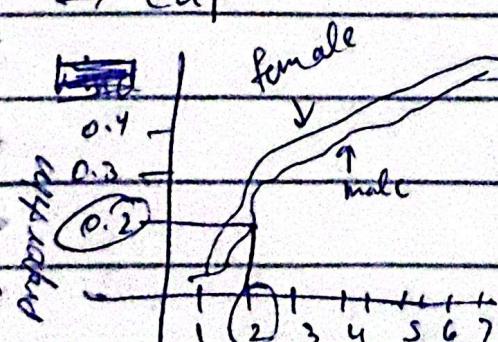
frequency

cdf

pdf

mdf

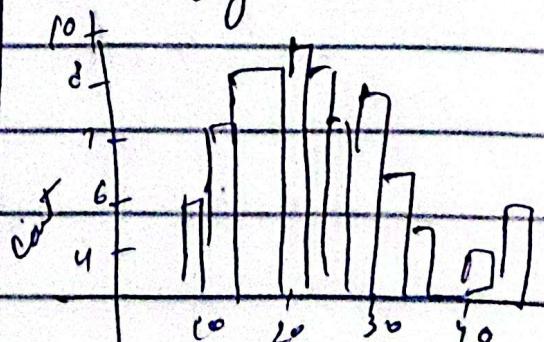
$\rightarrow$  cdf



20  $\rightarrow$  bill are paid

0.2 percent be female

day = Sat ...



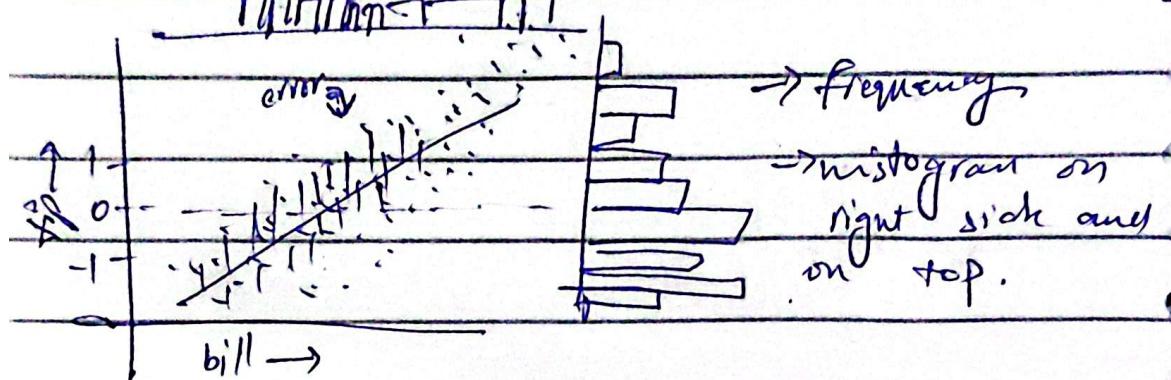
histogram  $\rightarrow$  several for every ~~and~~ day

Skewness

① Mean  $>$  Median  $>$  Mode  $\rightarrow$  Positive skew

② Mean  $<$  Median  $\rightarrow$  Negative skew

### B) Joint Plot



kind = "rplot", "scatter", "rcg", "hist", "resid"  
↓  
tells error distribution

→ Choose the graph on the basis on continuous or categorical data.

### Categorical Variable

1 → pick any number and find its  
2 frequency. If frequency > 1 then  
3 : it is categorical  
10 → can make group (grouping)

Wednesday

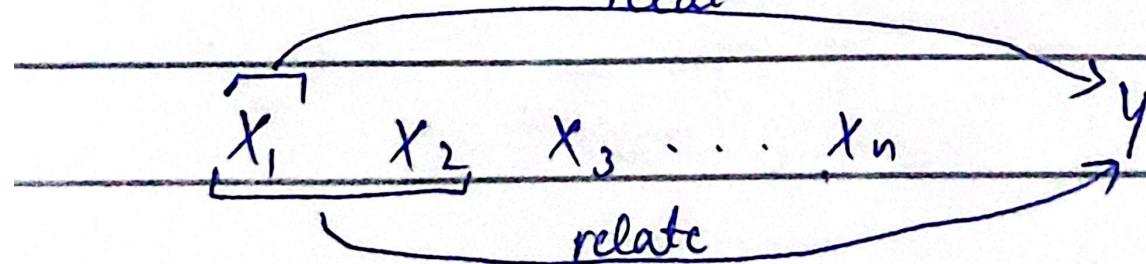
IDS

28-2

## Multivariate Data Analysis

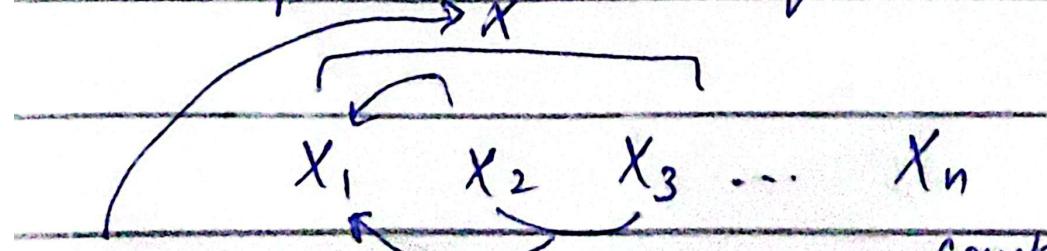
- Dependence Technique (Dependent & Independent)

relate



- Interdependence Technique

combine



- ① Factor analysis - Reduce the # of variables
- ② Cluster analysis - ~~combine~~ / Study together but not combined

- Scatter plot can be used to show multiple variables in 2D graph plane.

It uses continuous variables as x and

y and rest of any categorical variable

as: ① hue ② style ③ row/column or both

to show multiple variables on 2D graph.

① hue : means different color / shade

② style : (\*, o, •, △)

③ row/column : one graph, plan is divided

+	+
+	+

2 row x 2 col

Thursday

IDS

21-2

- Swarm plot
  - Playfair plot (9 plots)
    - r values
  - Graph Description (Explanation).
    - ↳ Explain graph with well-known words and used different words each time even for same meaning
- [Important for IELTS]

Wednesday.

IDS

6-3

Machine Learning - Automates analytical model

- prediction on previous / existing data

- learn from experience

Train the machine how to learn

Machine Learning

Structure Data

Data that can be represented in tabular form

Deep Learning

Unstructured Data

Data that cannot be represented in tabular form

- Given features from data

- Given data and extract features on their own

(Gedachten auf Fit)

R-Squared

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

$$r = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{\sqrt{n\Sigma x^2 - (\Sigma x)^2} \sqrt{n\Sigma y^2 - (\Sigma y)^2}}$$

$$\hat{y} = \frac{n\Sigma xy - \Sigma x \Sigma y}{n(\Sigma x^2) - (\Sigma x)^2}$$

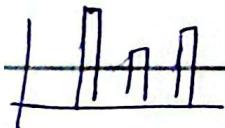
$$r = \frac{n\Sigma xy - \Sigma x \Sigma y}{\sqrt{[n\Sigma x^2 - (\Sigma x)^2][n\Sigma y^2 - (\Sigma y)^2]}}$$

$+1 \rightarrow$  positive relation  
 $-1 \rightarrow$  negative relation

$-1 \leq r \leq +1$

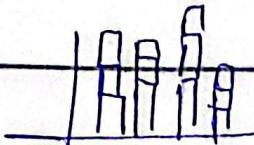
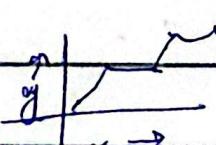
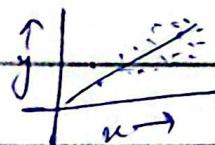
- Univariate Data / Variable (1 Variable)

- ① Bar chart    ② Pi chart    ③ Histogram



- Bivariate Data (2 variables)

- ① Scatter plot    ② Line plot    ③ Stacked Bar chart



- Multivariate Data (3 or more variables)

① style     $\circ \rightarrow \Delta \rightarrow *$

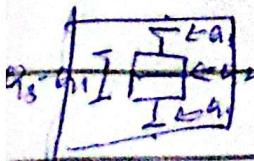
② hue, color, saturation

④ Row, columns increase

- Univariate Data Analysis

- ① Bar chart    ② Histogram    ③ Pie chart

- ① Box plot    ② Density Plot



## Bivariate Data Analysis

① Scatter Plot    ② Line Plot    ③ Count Plot

④ Box Plot

⑤ Swarm Plot



## Multivariate Data Analysis (more than 2 var)

① Interdependent Technique

② Interdependent technique → understand the structure

↳ Factor analysis → combine highly correlated features together

↳ Cluster analysis → combine similar features

Increasing trend ✓

Decreasing trend ✓

① Increased dramatically ✓

Slowed down noticeably ✓

② Tripled ✓

Modest — steady ✓

③ Increased strikingly ✓

Leveling off — steady ✓

④ Steep upward trajectory ✓

Declined ✓

⑤ Surged — increase suddenly ✓

Dwindled ✓

⑥ Climbing more slowly ✓

Plummeted ✓

⑦ Rising steadily ✓

Tapered off ✓

⑧ No sign of leveling off ✓

Decreased gradually ✓

~~Steeper~~ Fall down ✓

⑨ Increased exponentially ✓

Fell off ✓

⑩ Increased Rapidly ✓

Sank

⑪ Exponential Growth ✓

Waned

⑫ Accelerated ✓

To the Ground —

⑬ Rapid ascent ✓

Going downwards ✓

⑭ Substantial Increase ✓

Trend going down ~  
around straight slow ✓

## Graph code

ax.sns.displot(data=values, kind='ecdf') # hist, kde

→ hist: This values generate histogram plot.

↳ display the distribution of data through bin and provides insights into the f or pdf

→ kde: kernel estimation plot density estimation plot.

↳ used to estimate the pdf of a continuous var.

→ ecdf: empirical cumulative distribution function.

$$LR: \hat{y} = \theta_0 x + \theta_1, y = \beta_1 X + \beta_0 + \varepsilon \in \text{error}$$

Gradient Decent  
stop g' intersect

$$x_{n+1} = x_n + \alpha \Delta F(x_n)$$

$$F(x_n) = x^2$$

Cost Function (MSE)

$$x'_n = x + \alpha \Delta(x^2)$$

error<sup>2</sup> (given - actual)

$$x = x^2 + 2x$$

$$MSE = \frac{1}{n} \sum (y - \hat{y})^2$$

$$x = 2x + \alpha \Delta(2x)$$

$$y = mx + c$$

$$x = 2x + \alpha 2$$

$$MSE = \frac{\sum (y - \hat{y})^2}{n}$$

$$f(m, b) = \frac{1}{n} \sum (y_i - (mx_i + b))^2$$

$$\frac{\partial F}{\partial m} = \frac{1}{n} \sum -2(y - (mx + b))$$

$$\frac{\partial F}{\partial b} = \frac{1}{n} \sum -2(y - (mx + b))$$

$$P(x=x)$$

$$\frac{1}{10}, \frac{2}{10}, \dots, \frac{9}{10}$$

PMF (DDT)

0.1 0.1

0.2 0.3

PMF

1/20 1/20

PMF  $\rightarrow$  frequency (individually)

Discrete

$$P(x=x) f_{\text{total}}(a < x < b)$$

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

\*

PDF

cont  
rang

Normal distribution

$$= \frac{1}{6\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\sigma = \sqrt{\frac{\sum(x-\mu)^2}{N}}$$

$$\sigma = \sqrt{\frac{(32.74 - 29.52)^2 + \dots}{20}}$$

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



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

| sample mean - population mean

st /  $\sqrt{\text{number of sample}}$

$$z = \frac{1600 - 1570}{120/\sqrt{100}} = 2.51 > 1.96$$

$$b = \frac{n \sum xy - \sum x \sum y}{(n \sum x^2 - (\sum x)^2)}$$

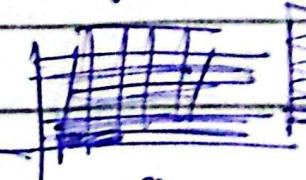
$$y = mx + b$$

$$a = \frac{\sum y - b \sum x}{n}$$

$$y = -4.4502 + 0.27x$$

$$R^2 = 1 - \frac{\sum (y - \hat{y})^2}{\sum (y - \bar{y})^2}$$

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{(n \sum x^2 - (\sum x)^2)(n \sum y^2 - (\sum y)^2)}}$$



$$NSB, \frac{\sum (y - \hat{y})^2}{n} \quad \hat{y} = mx + b$$

$$m = m - \alpha \frac{\partial d}{\partial m}$$

$$MSE = \frac{\sum (y - (mx+b))^2}{n}$$

$$b = b - \alpha \frac{\partial d}{\partial b}$$

$$MSE = \frac{\sum (y - (mx+b))^2}{n}$$

$$z\text{-score} = \frac{x-\mu}{\sigma} = \frac{\text{sample} - \text{pop}}{\text{st}}$$

$$MSE = \frac{\sum (y - mx - b)^2}{n}$$

$$z\text{-test} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$a = \frac{n \sum y - b \sum x}{n}$$

$$r = \frac{\sum xy - \sum x \sum y}{\sqrt{\sum x^2 - (\sum x)^2} \sqrt{\sum y^2 - (\sum y)^2}}$$

Wednesday

IDS

20-3

kNN → Supervised learning

↳ vote from majority      (most similar)

↳ used for classification      similarity

↳ regression

## Binary classification $\rightarrow$ better

- non-parametric algorithm - don't assume about data
- lazy learner - nearly no learning instead it store the dataset
- Find the distance of all of the points in the plane (1000, 1000) and then sort the distance points. Then check the first  $k$  points and then predict the classes.  $k = \text{any odd number}$ .

### Distance functions

① Euclidean

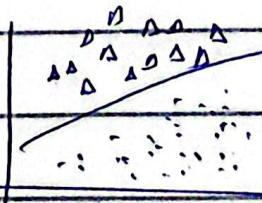
$$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

② Manhattan

③

Example

Linear classifier



N.P (3, 7),  $k = 3$

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$d =$$

③ A  
Sort = {3, 3, 6, 9, 5}  
 $\xrightarrow{k}$

x	y	label	d
7	7	A	4
7	4	A	5
3	4	B	3
1	4	B	3.6

Class = B

### Value of $k$

- Value of  $k$  must be odd (3, 5, 7, 11)
- Small values of  $k$  can make noise
- Larger values of  $k$

## Table

Expected <small>Input</small>	Actual
value[i] = -999	
value = 1000, i=101	

Friday

IDS

22-3

K-Means

- Clustering algo

- Intra cluster /within cluster

distant should be minimum

$n \rightarrow$  no. of cluster  $\{k=2, 3, \dots\}$

$\hookrightarrow$  no. of centroids = center points

$\rightarrow$  find the distance of cluster

of every point to centroid and the

point which have minimum distance with

centroid, group them e.g.  $|x_1 - c_1, x_2 - c_2|$

$\rightarrow$  find the mean distance within centroid

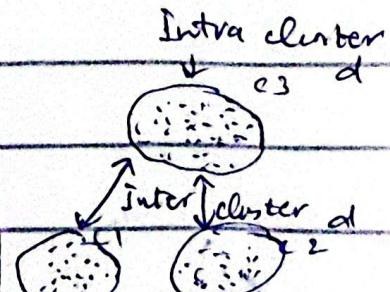
update the value

and move the centroid point in the middle

of the cluster

$\rightarrow$  Do this until there is some values of

average



$x \rightarrow$

$c \rightarrow$  cluster

→ ~~randomly~~  
initialize centroid  
values

no. of cluster

choosing the  $k$  value

- Elbow method — WCSS

Intra cluster  $\downarrow$

- Expectation Maximization

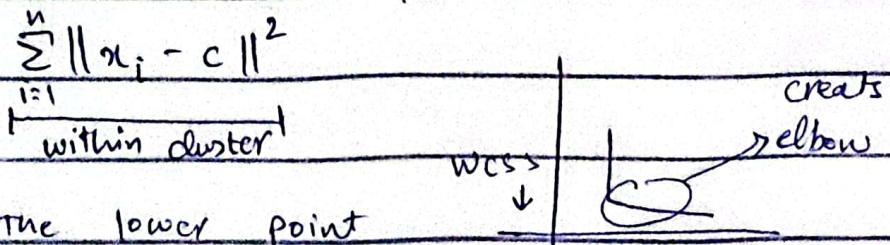
Inter cluster  $\uparrow$

$$d = \sum_{j=1}^c \sum_{i=1}^n \|x_i - c_j\|^2$$

data points                            clusters

- Elbow method — find the value of  $k$

$\hookrightarrow$  within cluster sum of squares (WCSS)



- choose the lower point

of  $(k)$   $\rightarrow$  optimal number of clusters  $\boxed{k \rightarrow \text{integer}}$

- Silhouette Method — used for verification mainly

$b \rightarrow$  distance of the point from cluster it does not belong

$a \rightarrow$  distance of the point from cluster it belongs  $\Rightarrow -1 < \text{value} < 1$

Kmean. inertia } return sum of squares

+ve  $\rightarrow$  good, -ve  $\rightarrow$  bad, 0  $\rightarrow$  overlapping

$\downarrow$  optimal

$\downarrow$  not optimal

Tuesday

IDS

27-3

## Evaluation Metrics

Dataset

train-test-split()  
Feed to model

x-train  
y-train  
x-test  
y-test

Training → x-train  
y-train

Testing → x-test (y-predict) → compare y-test → confusion matrix

## Classification → Confusion Matrix

		Predicted			
		P: Yes	N: No	Type 1 error	Type 2 error
Actual	T	TP	FP	False Positive	True Positive
	F	FN	TN	False Alarm	True Negative
				Type 1 error	Type 2 error

① Accuracy =  $\frac{\text{Correct}}{\text{Total}}$

$\uparrow$  Type 1 error  $\rightarrow$  Overfit  
 $\uparrow$  Type 2 error  $\rightarrow$  Underfit

Example		Actual	Dog	TP : 3	TN : 4
Predicted				FN : 2	FP : 2
3	2				
1	4				

$$\text{Accuracy} = \frac{3+4}{10} = \frac{7}{10} = 70\%$$

## Classification → Confusion Matrix

Predicted				Type 1 error	False Positive	False Alarm
P: Yes	N: No	T	F			
Actual T	TP	FP	F	Type 1 error	False Positive	False Alarm
F	FN	FN	T	Type 2 error	↑ Type 1 → Overfit	↑ Type 2 → Underfit

(1) Accuracy =  $\frac{\text{Correct}}{\text{Total}}$

Example		Actual	Dog	TP : 3	TN : 4
Predicted		3	2	FN : 2	FP : 2
		1	4		

$$\text{Accuracy} = \frac{3+4}{10} = \frac{7}{10} = 70\%$$

Friday

## IDS

2893

Naïve Bayes - Bayes Theorem

- Conditional Probability

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

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$   
 $A = X \quad B = y \quad X \quad y$   
 $P(A|B) \quad P(B|A) \cdot P(A) \quad P(B)$

① Independent

② Equal Imp

$$\prod_{i=1}^n$$

Posterior = Prior × likelihood

↑  
Product

$X = X_1, X_2, X_3, \dots, X_n$  evidence

$$P(y|X_1, X_2, X_3, \dots, X_n) = \frac{P(X_1, X_2, X_3, \dots, X_n|y) \cdot P(y)}{P(X_1, X_2, X_3, \dots, X_n)}$$

- Zero Frequency problem  $\Rightarrow$  Laplacian Smoothing

↳ If data not exist of a feature then

the probability is zero which makes the whole probability zero

NB types

① Multinomial    ② Bernoulli    ③ Gaussian

Ex#1

$$P(\text{write} \mid \text{Infected}) = \frac{\text{Total}}{7}$$

$$P(\text{write} \mid \text{Clean}) = \frac{P(\text{Clean})}{7} = \frac{3}{7}$$

$$P(\text{Executable} \mid \text{Infected}) =$$

$$P(\text{Non-exe} \mid \text{Infected}) = \frac{P(\text{Infected})}{4}$$

	write	executable	large
P(Clean)	2/3	1/3	2/3 $\Rightarrow$ 3/4
P(Infected)	1/4	4/4 = 1	<del>1/4</del> $\Rightarrow$ 1/5

$$P(\text{Clean} \mid \text{Yes}) = \frac{2}{3} \times \frac{1}{3} \times \frac{3}{4} = 0.071 \uparrow$$

$$P(\text{Infected} \mid \text{No})$$

$$P(\text{Clean} \mid \text{No}) = \frac{1}{4} \times \frac{1}{4} \times \frac{3}{7} = 0.021$$

Wednesday

IS

3-4

Wednesday

IDS

3-4

### Spam Filter (Important)

- Types : Bayesian Filter

$$P(\text{spam}|\text{word})$$

$$\Pr = \frac{\Pr(\text{word}|\text{spam}) \Pr(\text{spam})}{\Pr(\text{word})}$$

$$\text{Total} = [p_1 * p_2 * p_3]$$

Ham: send me review

$$P(S/W) = P(W/S) * P(S)$$

Spam: Send me password

$$P(W/S \cup H)$$

$$P(W/S \cup H) = P(W/S) * P(S) + P(W/H) * P(H)$$

$$P = P_1 P_2 \dots P_n$$

$$P_1 P_2 P_3 \dots P_n + (1 - P_1)(1 - P_2) \dots (1 - P_n)$$

Laplace Smoothing (+2 in classes) + 1 in features  
every word is a feature

$E_n = \text{"Review us now"}$

$$P(S/E_n) = \frac{P(E_n/S) * P(S)}{P(E_n)}$$

$P(E_n/S)$  zero for word which is not in  $E_n$  ( $1 - P$ )

$\Pr(\{1, 0, 1, 0, 0, 0\} | \text{spam})$

one on the  
for the word which exists in  $E_n$

( $p$ )

- Threshold / Tolerance = 5%

Friday

SQE

5-4

Quiz - 2

Friday

IDS

5-4

### Data Preprocessing

- Raw Data  $\downarrow$  apply  $\uparrow$  (80% Time)

Prepared data  $\xrightarrow{\text{feed}}$  algorithm

### Data Cleaning

aka cleansing, scrubbing

- remove and fix data

Raw Data

$\downarrow$  Data Preprocessing

Prepared Data

$\downarrow$  Feed

Algorithm

- remove duplicate data

} affect models

- irrelevant data

} only kNN can work  
on empty cells

① variance is zero

$\text{NaN} \rightarrow$  Empty cells

② gives no information

e.g. email, phone

Note:- If all of the data is 100% same

\* different then remove that column

Python → [NaN] → empty cells

## Missing Values

① Missing Completely at Random (MCAR)

↳ missing data cannot be linked/related with other features

② Missing At Random (MAR)

↳ missing data can be related with other features

③ Missing Not At Random (MNAR)

↳ Unobserved data, cannot find any relation

## Handling Missing Value

① Deleting the missing values

- If 2% of 100k entries are missing then deleting is a good approach. If MAR and MCAR then delete it.

② Imputing the missing values

- Replace with arbitrary value

- Replace with mean

- Replace with Mode (for categorical data)

- Replace with Median (for outliers)

- Replace with Previous fill, Replace with next word

Data Transformation

[String → Number] fill

① Label ② One Hot Encoding

var types  
① Nominal (Categorical) ② Ordinal

~~But~~ Categorical Data

## One Hot Encoding (Example)

Color	-red	-green	-blue
Red	1	0	0
Green	0	1	0
Blue	0	0	1

Vectorize form  
One Hot encoding vector

Text → Number → NLP  
→ Counter Vectorizer

## Feature Scaling

### ① Min-Max Normalization

Salary	Age
--------	-----

10k	22
-----	----

20k	27
-----	----

50k	39
-----	----

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

### ② Standardization

Very big values affect the model

small values neglected by model

$$x' = \frac{x - \bar{x}}{s}$$

[must]

Normalization vs Standardization

no Gaussian Distribution

Follow Gaussian Dist.

Friday

IDS

194

## Data Preprocessing

### - Feature Generation & Selection

#### • Selection

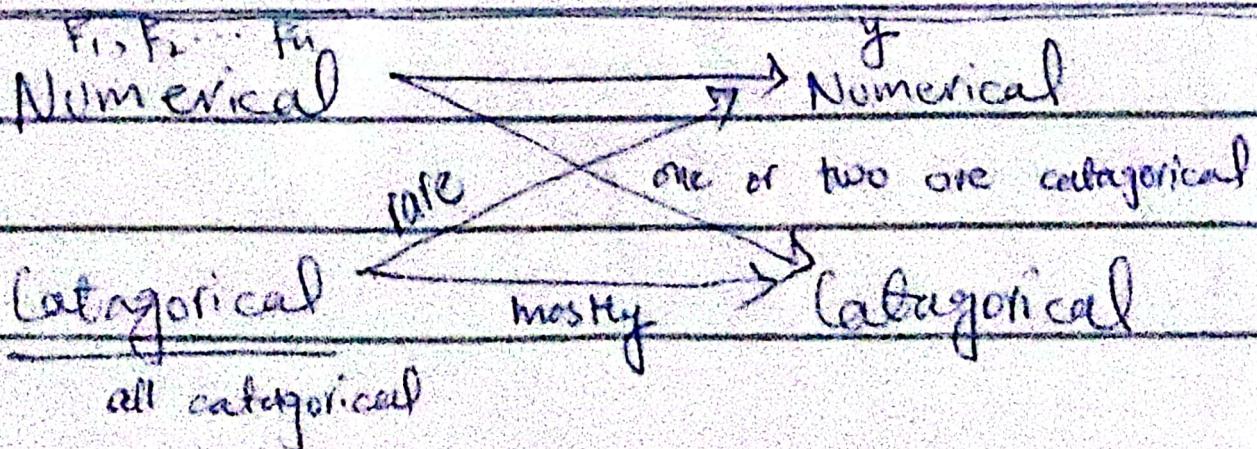
##### ① Unsupervised

- Remove with zero or constant variance
- Remove missing values column.
- High multicollinearity  $\Rightarrow$  high correlation cause redundancy. - Remove them. [less variety]

##### ② Supervised

###### ① Filters Method aka Single Factor Analysis

- ↳ correlate ~~the~~ each feature with target / label



Not always increases accuracy

## ① Numerical $\rightarrow$ Numerical (Regression)

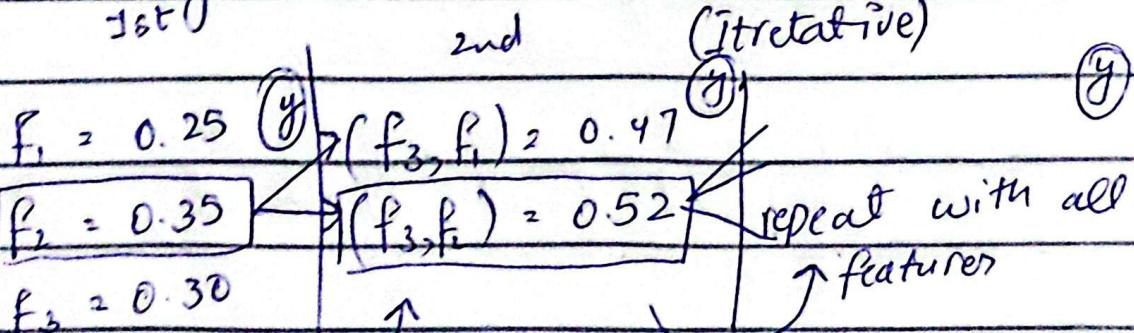
- Find correlation

## ② Numerical $\rightarrow$ Categorical (Classification)

- $\hookrightarrow$  ① Information Gain / Mutual Information

## ② Wrappers

- Greedy Search



- Backward Elimination

- $\hookrightarrow$  can go backward and eliminate selected features

## ③ Embedded Method

- Decision Tree

Regression

Classification

Information Gain

Regression D.T  $\rightarrow$  Classification

Information Gain

Selection Model  
Algo.

Learning Model  
Algo.

Information Gain → ASM → Attribute Selection Method

Decision Tree

Information Gain  
Gini Index  
Gain Ratio

① Node →

Information Gain

② Edge →

Depends  
Entropy

③ Leaf →

Entropy → Randomness

Calculate Entropy for 1 → Yes only = 0

$$E(S) = \sum_{i=1}^n -p_i \log_2 p_i$$

$\log_2 \rightarrow \log$

of base] - Find target var classes  
 Imp! 1 2 e.g. Yes = 9 } 14  
 9, 5 No = 5 }

$$E(5, 9) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

Yes                  No

Entropy for multiple entropy attributes

$$P_1, P_2, P_3, \dots, P_n \rightarrow Y$$

entropy

$$E(T, X) = \sum_{c \in X} P(c) E(c)$$

$$F(a, b) = \left( -\frac{a}{a+b} \log_2 \frac{a}{a+b} \right) + \left( -\frac{b}{a+b} \log_2 \frac{b}{a+b} \right)$$

-ve in Entropy:

- $\log_2$  of ratio is always -ve so we put +ve in it

Random Forest  $\leftarrow$  DT<sub>1</sub>, DT<sub>2</sub>, ..., DT<sub>N</sub>

Information Gain  $\rightarrow$  IG<sub>T</sub>

$\Rightarrow$  Information Gain (T, X) = Entropy(T) - Entropy(X)

T  $\rightarrow$  Total attri

X  $\rightarrow$  single attri

$\Rightarrow$  Information Gain = Entropy(before) -  $\sum_{j=1}^k$  Entropy(j, after)

- Select with the maximum IG<sub>T</sub>

Wednesday

IDS

24-4

Data Transformation

↳ Dimensionality Reduction

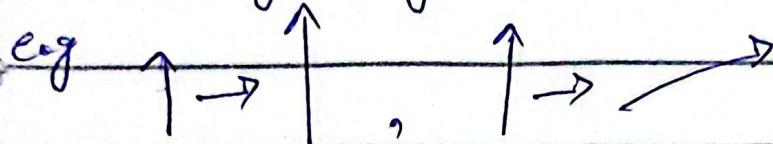
⇒ Transformation → change of ~~the~~ shape

## ① Linear Transformation (Geometry)

- straight line transformation

- size may change

- no curve, no bend



## Eigenvectors and Eigenvalues (Linear algebra)

Matrix =  $A$

Vector =  $\vec{v}$  single column

$$\begin{bmatrix} 2 & 3 \\ 4 & 5 \end{bmatrix} * \begin{bmatrix} 4 \\ 9 \end{bmatrix} = \begin{bmatrix} 20 \\ 45 \end{bmatrix}$$

Dataset

scaled

always vector

- Resultant vector is either scale-up or scale-down without changing direction linearity.
- Scale-up / Scale-down with an integer and that integer is called eigen value.

→ Eigenvector

$$A \cdot \vec{v} = \lambda \cdot \vec{v}$$

↑ scalar

find  $\vec{v}$ ?

$\lambda \rightarrow$  eigenvalue

$v \rightarrow$  eigenvector (Scaled version)

$$A\vec{v} - \lambda\vec{v} = 0$$

$I \rightarrow$  Identity matrix

$$\vec{v} (A - \lambda I) = 0$$

$\lambda \rightarrow$  can have more than

$$\vec{v} \neq 0 \text{ then } A - \lambda I = 0$$

one values

Use in DS

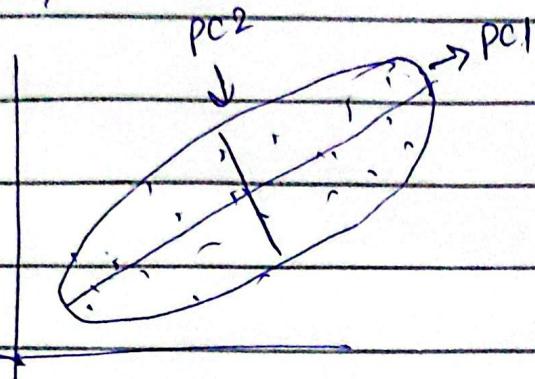
→ Determine a set of common variables

# Dimensionality Reduction

## ① Principal Component Analysis

### ② Singular Value Decomposition

①  $\rightarrow$



$$A - \lambda I = 0$$

$$\lambda = \lambda_1, \lambda_2$$

$$PC_1 > PC_2$$

$v_1$        $v_2$

Orthogonal  
90°

Variance

→ each feature  
are independent

Principal Component

→ each depends on  
each other (combine)

Friday

SQE

26-4

JMeter - Quiz on Friday

Friday

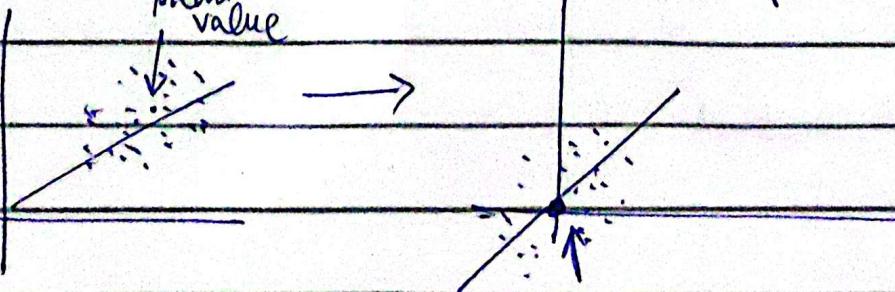
IDS

26-4

PCA  $\Rightarrow$  also calculates co-variance rather

than variance of separate features only

Mean value



→ Move mean value

Mean value

to center (0,0) for 2-Distribution

→ v) cannot apply PCA directly / standardization / without scaling

## → Mathematically :

Example : Apply PCA to find 2-3 important out of 4.

i) find mean & std of each feature col.

ii) find  $x_{\text{new}} = \frac{x - \mu}{\sigma}$

iii) find variance

$$\sigma^2 = \frac{\sum (x - \mu)^2}{N}$$

iv) find co-variance

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

v) calculate cov matrix  $\text{cov}_{ij} = \text{cov}_{ji}$   $n = \text{no. of obsr.}$

$f_1$	$f_2$	$f_3$	$f_4$	
$f_1$	$V_{f_1}$	cov	cov	cov
$f_2$	cov	$V_{f_2}$	cov	cov
$f_3$	cov	cov	$V_{f_3}$	cov
$f_4$	cov	cov	cov	$V_{f_4}$

vi) Calculate eigenvalues and eigenvectors

$$\begin{bmatrix} \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix} = 0$$

right vector

standardized form of features

Value off  $\lambda = \text{no. of columns}$

0. --

Biengal

Feature matrix \* eigen. vector = Transformed Data  
no. of  $k$

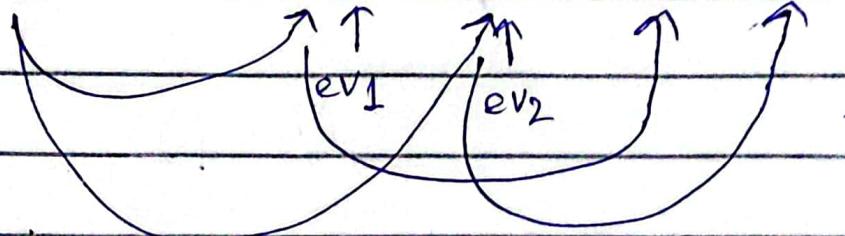
eig  $\downarrow$

$k=2$

$nf_1$

$nf_2$

$$\begin{bmatrix} - & - & - & - \\ - & - & - & - \\ - & - & - & - \\ - & - & - & - \end{bmatrix} \begin{bmatrix} - & - & - \\ - & - & - \\ - & - & - \\ - & - & - \end{bmatrix} \begin{bmatrix} - & - & - \\ - & - & - \\ - & - & - \\ - & - & - \end{bmatrix} \begin{bmatrix} - & - & - \\ - & - & - \\ - & - & - \\ - & - & - \end{bmatrix}$$



## Importance

- Improve accuracy
- Less interpretable
- Standardization is necessary

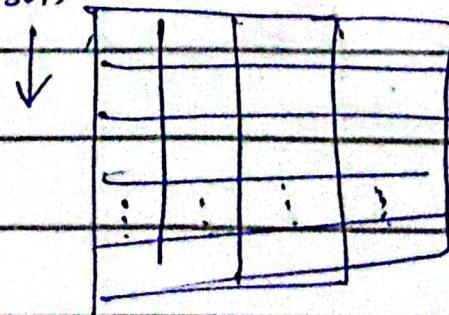
## Recommendation System

## Matrix Factorization

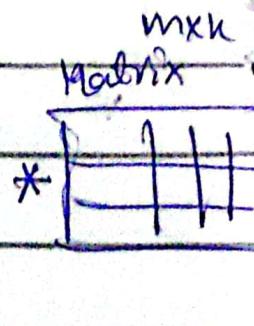
$$[5 \times 5] \rightarrow [5 \times 2] * [2 \times 5]$$

Users Movie

$K \ll n, m$



User



## Singular Value Decomposition (SVD)

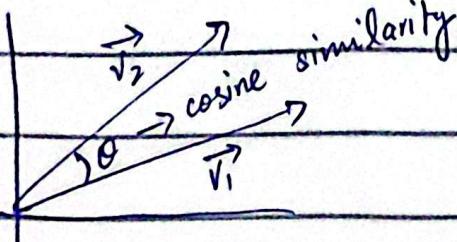
$$R = P \Sigma Q^T$$

$$\left[ \begin{array}{c} \\ \\ \end{array} \right] = \left[ \begin{array}{c} \\ \\ \end{array} \right] * \left[ \begin{array}{cc} 0 & \theta^\circ \\ & 1 \end{array} \right] *$$
  
$$\left[ \begin{array}{c} \\ \\ \end{array} \right]$$

## Similarity Measure

Object 1      Object 2  
vectors      vectors

### - Cosine Similarity



- similarity between 2 vectors in direction
- the angle decides the similarity between vectors

### - Convert features into vectors

### - Dot Product the vectors

$$\text{similarity}(A, B) = \cos \theta = \frac{A \cdot B}{|A||B|}$$

very similar      not similar

Doc 1 : \_\_\_\_\_  $\Rightarrow \vec{v}_1 = [1 \ 1 \ 0 \ 1 \dots]$

Doc 2 : \_\_\_\_\_  $\Rightarrow \vec{v}_2 = [1 \ 0 \ 1 \dots]$

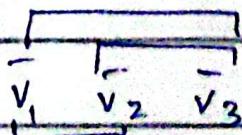
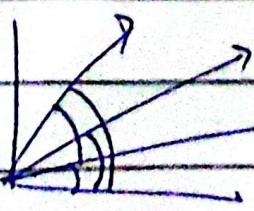
→ BOW returns the unique words in Doc.

→ words appearing in the doc are shown

only one time but their frequency should

be in vectorize form

- For 3 and more



Users	Movies		
U <sub>1</sub>	M <sub>1</sub>	M <sub>2</sub>	M <sub>3</sub>
U <sub>2</sub>	2	3	5
U <sub>3</sub>	1	2	5

→ For user

↑

For movie

- Bias

Predict — Actual  
Gap

- Variance = Predicted ] Spread  
Values

### • Overfitting

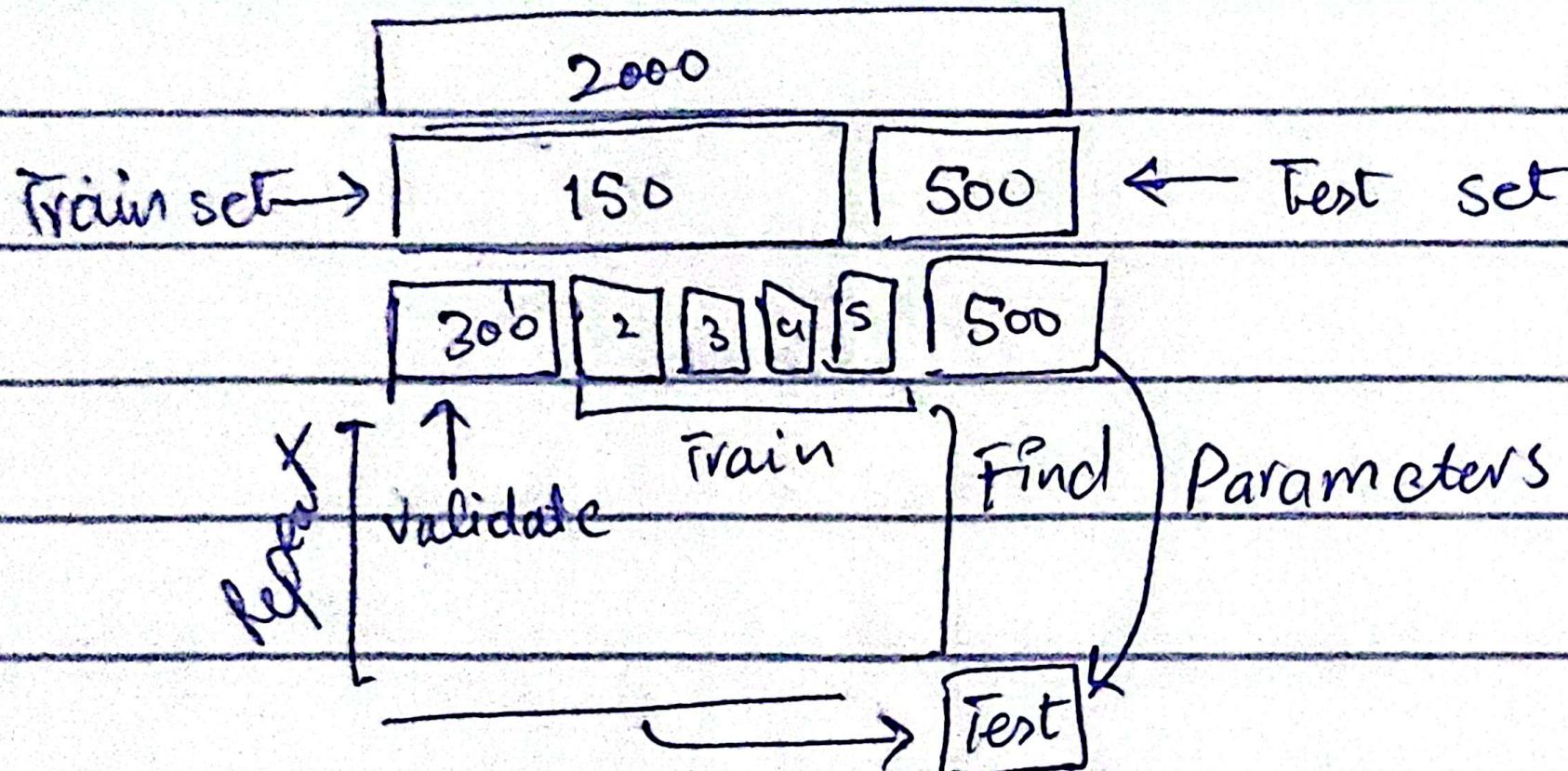
- Model failed to generalize, instead it specifies the features.
- Model memorize the features.
- Model performed well on training data but performed poor in train-test data.

### • Underfitting

- Model doesn't learn in any way, no feature learnt.
- Model unable to capture features.
- Model performed poorly on both training and testing dataset.

## 10-Fold

### k-Fold Cross-Validation



## Confusion matrix of multi-class classification

- No. of classes x No. of classes

Actual

		A	B	C	
Predicted	A	TP	FP		Sum
	B		TN		
	C	FN			Sum

sum  
sum

Actual

	TP	FP = Sum	
Predicted	FN	TN = Sum	$\Rightarrow$ 2x2 matrix
	Sum		

$\rightarrow$  Multi-class

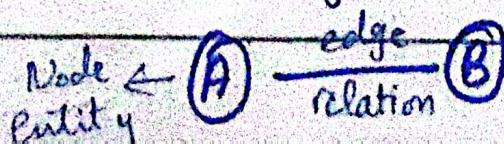
- < Precision, Recall

TN	FN	TN	> Accuracy
FP	TP	FP	$\rightarrow$ Binary-classification
TN	FN	TN	< Accuracy

## Graph Data Science

- data presented as graphs

- Social networking e.g.



# Centrality Measures in Social Network Analysis

## ① Degree centrality

- one hop / jump      - find very connected individual

## ② Betweenness Centrality

$\text{Gst}(v) \leftarrow \text{each pair}$   
 $\text{Gst}$

## ③ Closeness Centrality

$$CC(i) = \frac{N-1}{\sum_{j=1}^{N-1} d(i,j)}$$

	A	B	C	farmer	$\frac{N-1}{\sum_{j=1}^{N-1} d(i,j)}$
A	0			→ sum	$\frac{(3-1)}{\text{Sum}}$
B		0		→ sum	$\frac{(3-1)}{\text{Sum}}$
C			0	→ sum	$\frac{(3-1)}{\text{Sum}}$