

$$y_1 = \begin{bmatrix} 0.91 & 0.39 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \end{bmatrix} =$$

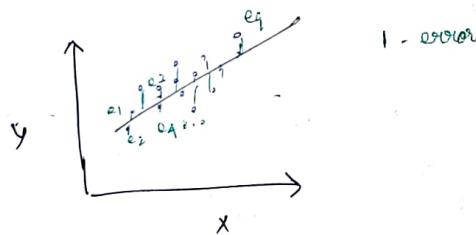
16/3

## Regression Analysis

Primary objective : to predict data

Linear regression

$x \rightarrow$  independent variable  
 $y \rightarrow$  dependent variable



$$\text{residue} = e_1 + e_2 + \dots + e_q$$

$$y = a_0 + a_1 x + e$$

$a_0 \rightarrow$  Intercept

$a_1 \rightarrow$  Slope

$e \rightarrow$  error term

$$y_i = a_0 + a_1 x_i + e_i$$

$$\Rightarrow e_i = y_i - (a_0 + a_1 x_i)$$

$$\text{Sum of errors squared, } E = \sum_{i=1}^n e_i^2$$

$$= \sum_{i=1}^n [y_i - (a_0 + a_1 x_i)]^2$$

Sum of errors squared,

$$= \sum_{i=1}^n e_i^2$$

$$= \sum_{i=1}^n [y_i - (a_0 + a_1 x_i)]^2$$

Minimisation fn.

$$J(a_1, a_0) = \sum_{i=1}^n [y_i - (a_0 + a_1 x_i)]^2$$

$$a_1 = \frac{\bar{xy} - \bar{x} \cdot \bar{y}}{(\bar{x}^2) - (\bar{x})^2}$$

$$a_0 = \bar{y} - a_1 \bar{x}$$

Q.	$x_i$ (Week)	$y_i$ (Sales in thousands)	$xy$	$x^2$
	1	1.2	1.2	1
	2	1.8	3.6	4
	3	2.6	7.8	9
	4	3.2	12.8	16
	5	3.8	19.0	25

$$\bar{x} = 3.3 \quad \bar{y} = 2.52 \quad \bar{xy} = 12.44$$

$$8.88$$

$$\bar{x}^2 = 11$$

$$a_1 = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - \bar{x}^2} = \frac{8.88 - 3(2.52)}{11 - 9} \\ = \frac{1.32}{2} = 0.66$$

$$a_0 = \bar{y} - a_1 \bar{x} = 2.52 - 0.66(3) \\ = 0.54$$

$$\Rightarrow y = 0.54 + 0.66x$$

$\Rightarrow$  Salary for 7<sup>th</sup> week

$$= 0.54 + 0.66(7)$$

$$y = 5.16$$

$$\text{for 12<sup>th</sup> week } y = 8.46$$

Matrix Representation

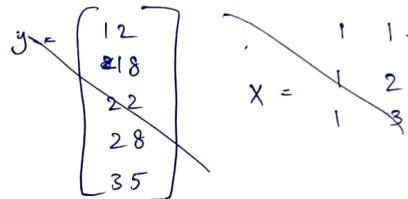
$$Y = X\alpha + e$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \end{bmatrix} + \begin{bmatrix} e_0 \\ e_1 \\ \vdots \\ e_n \end{bmatrix}$$

$$\Rightarrow \alpha = ((x^T x)^{-1} x^T) Y$$

Q. Consider . . . no. of hours spent in a library. Predict the no. of hours that would be spent in the 7<sup>th</sup> and 9<sup>th</sup> week.

Week	Hours Spent
1	12
2	18
3	22
4	28
5	35



x <sub>0</sub>	y	x <sup>2</sup>	x y
1	12	1	12
2	18	4	36
3	22	9	66
4	28	16	112
5	35	25	175

$$\bar{x} = 3, \quad \bar{y} = 23 \quad \bar{x^2} = 11 \quad \bar{xy} = 80.2$$

$$a_1 = \frac{\bar{xy} - \bar{x}\bar{y}}{\bar{x^2} - \bar{x}^2} = \frac{80.2 - 3(23)}{11 - 3^2}$$

$$= \frac{11 \cdot 2}{2} = 5.6 ; \quad a_0 = \bar{y} - a_1 \bar{x} = 6 \cdot 2 \\ \Rightarrow y = 6 \cdot 2 + 5.6x$$

$$\text{For 7<sup>th</sup> week, } y = 45.4$$

$$\text{9<sup>th</sup> week, } y = 56.6 \text{ hours}$$

21/3 Multiple Linear Regression

$$y = f(x_1, x_2, \dots, x_n)$$

$$= a_0 + a_1 x_1 + a_2 x_2 + \dots + a_n x_n + \epsilon$$

$$\hat{a} = ((x^T x)^{-1} \cdot x^T) y$$

	$x_1$	$x_2$	$y$
1	2	4	6
2	5	8	12
3	8	12	20
4	12	-	-

$$x = \begin{bmatrix} 1 & 1 & 4 \\ 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 12 \end{bmatrix}, x^T = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 4 \\ 1 & 4 & 12 \end{bmatrix}$$

$$x^T x = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 4 \\ 1 & 4 & 12 \end{bmatrix} \begin{bmatrix} 1 & 1 & 4 \\ 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 12 \end{bmatrix} = \begin{bmatrix} 10 & 10 & 29 \\ 10 & 30 & 86 \\ 10 & 86 & 249 \end{bmatrix}$$

$$(x^T x)^{-1} = \begin{bmatrix} 3.15 & -0.59 & -0.3 \\ -0.59 & 0.2 & 0.016 \\ -0.3 & 0.016 & 0.054 \end{bmatrix}$$

$$(x^T x)^{-1} = \begin{bmatrix} 3.15 & -0.59 & -0.3 \\ -0.59 & 0.2 & 0.016 \\ -0.3 & 0.016 & 0.054 \end{bmatrix}$$

$$(x^T x)^{-1} x^T = \begin{bmatrix} 3.15 & -0.59 & -0.3 \\ -0.59 & 0.2 & 0.016 \\ -0.3 & 0.016 & 0.054 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 4 \end{bmatrix} = \begin{bmatrix} 1.36 & 0.47 & 1.02 \\ 0.32 & 0.1 & 0.155 \\ -0.065 & 0.002 & 1.85 \end{bmatrix}$$

$$= \begin{bmatrix} 1.36 & 0.47 & 1.02 & -2.81 \\ 0.32 & 0.1 & 0.155 & 0.402 \\ -0.065 & 0.002 & 1.85 & -0.128 \end{bmatrix}$$

$$(x^T x)^{-1} x^T y = \begin{bmatrix} 1.36 & 0.47 & 1.02 & -2.81 \\ 0.32 & 0.1 & 0.155 & 0.402 \\ -0.065 & 0.002 & 1.85 & -0.128 \end{bmatrix} \begin{bmatrix} 1 \\ 6 \\ 8 \\ 2 \end{bmatrix} = \begin{bmatrix} 6.72 \\ 2.964 \\ 47.239 \end{bmatrix}$$

$$\hat{a} = \begin{bmatrix} 6.72 \\ 2.964 \\ 47.239 \end{bmatrix} \quad \begin{bmatrix} -1.69 \\ 3.48 \\ -0.05 \end{bmatrix}$$

$$y = 1.69 + 3.48x_1 - 0.05x_2$$

Q. Using multiple regression analysis

x

z = Equity

x = net sales

y = Asset

z is dependent variable

x, y are independent variables

## Polynomial Regression

$$y = ax^b$$

→ Non-linear

$$\log y = \log(a x^b)$$

$$\Rightarrow \log y = \log a + b \log x$$

$$y = a_0 + a_1 x + a_2 x^2 + \dots$$

$$\begin{matrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{matrix} \quad \text{Matrix format of polynomial reg.}$$

$$\begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \\ \sum x_i^2 y_i \end{bmatrix}$$

$$a = X^{-1} B$$

	$x_i$	$y_i$	$\frac{\sum x_i y_i}{n}$	$\frac{x_i^2}{n}$	$\frac{x_i^3}{n}$	$\frac{x_i^4}{n}$	$\frac{\sum x_i^2 y_i}{n}$
1	12	8					
2	18	12					
3	22	16					
4	28	36					
	35	42					

$$x^T x = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 12 & 18 & 22 & 28 & 35 \\ 8 & 12 & 16 & 36 & 42 \end{pmatrix} \begin{pmatrix} 1 & 12 & 8 \\ 1 & 18 & 12 \\ 1 & 22 & 16 \\ 1 & 28 & 36 \\ 1 & 35 & 42 \end{pmatrix} = 5 \times 5 \quad 5 \times 3$$

$$\begin{pmatrix} 5 & 115 & 114 \\ 115 & 2961 & 8142 \\ 114 & 3142 & 3524 \end{pmatrix}$$

$$\Rightarrow a = \begin{pmatrix} 1 & 10 & 30 \\ 10 & 30 & 100 \\ 30 & 100 & 300 \end{pmatrix}^{-1} \begin{pmatrix} 29 \\ 96 \\ 338 \end{pmatrix} \Rightarrow a = \begin{pmatrix} -0.75 \\ 0.95 \\ 0.75 \end{pmatrix}$$

$$y = -0.75 + 0.95x^2 + 0.75x^4$$

## Logistic Regression

$$P(x) = a_0 + a_1 x$$

$$y = \frac{1}{1 + \exp(-a_0 - a_1 x)}$$

Q. Let's assume a binomial logistic regression where the classes are pass and fail. The student dataset has entrance based on historic data of those who are selected & not selected. Based on logistic regression values of learned parameters,

$$a_0 = 1, a_1 = 8 \quad \text{assuming marks of } n=60,$$

compute the resultant values

$$a_0 + a_1 x = 1 + 8(60) = 481$$

$$y_{\text{Pass}} = \frac{1}{1 + \exp(-481)} = 1 \xrightarrow{0.44?}$$

> 0.5 selected

else not selected

## Validation of regression method

### ① Mean Absolute Error (MAE)

$$\frac{1}{n} \sum_{i=0}^{n-1} |y_i - \hat{y}_i|$$

$y_i$  = Actual target,  $\hat{y}_i$  = Predicted output

### ② Mean Squared Error (MSE)

$$\Rightarrow \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2$$

### ③ Root Mean Squared Error (RMSE)

$$= \sqrt{\text{MSE}}$$

### ④ Relative MSE

$$\frac{\sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2}$$

C. coefficient of variation,  $CV = \frac{RMSE}{\bar{Y}}$

Standard error:  $y_i - \hat{y}$

Q.  $x_i$   $y_i$  (Actual Sales)

i1	80
i2	90
i3	100
i4	110
i5	120

Consider 2 fresh items i6 & i7 whose actual values are 80 & 75 respectively. Regression model predicts the values of i6 & i7 as 75 & 85 respectively.

Find MAE, MSE, RMSE, Relative MSE, CV

$n=2$

$$MPE = \frac{1}{2} (175 - 80) + (85 - 75)$$

$$= \frac{1}{2} (5 + 10) = 7.5$$

$$MSE = \frac{25 + 100}{2} = \frac{125}{2} = 62.5$$

$$RMSE = \sqrt{62.5} = 7.9$$

$$\text{Relative MJE} = \frac{\sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2}$$

$$= \frac{(80 - 75)^2 + (75 - 85)^2}{(80 - 77.5)^2 + (75 - 77.5)^2}$$

$$= \frac{25 + 100}{2.25 + 3.25}$$
~~$$= \frac{125}{5.5} = 22.72$$~~

$$= 10$$

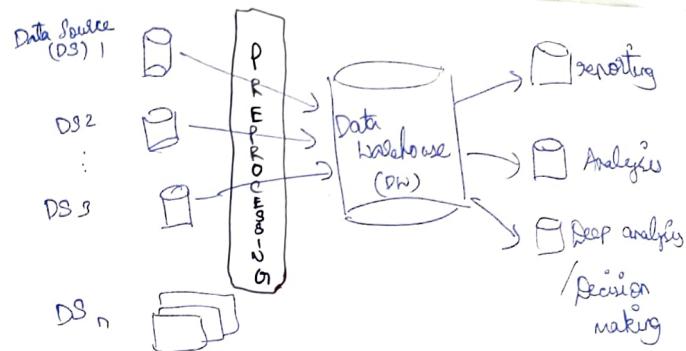
$$CV = \frac{RMSE}{\bar{y}} = \frac{7.9}{77.5} = 0.102$$

23) 3

### Data Warehouse (DW)

#### Characteristics

1. Subject-oriented (Sale/Marketing/Financial etc.)
2. Integrated
3. Non-volatile



### Online Analytical Processing - OLAP

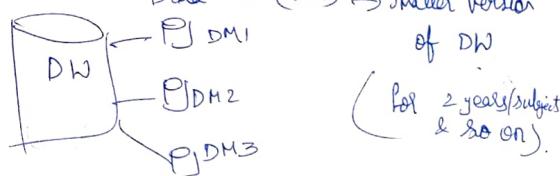
#### 1. Time variant (Date / Month / Year)

Date warehouse can contain

- metadata
- summary of data

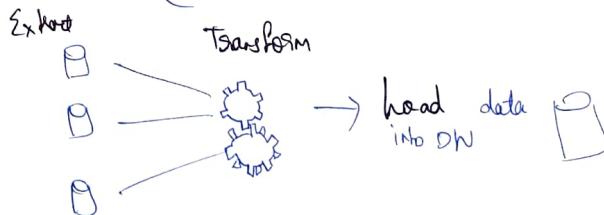
→ Tools are used in DW for analysis & reporting  
(Business Intelligence)

Data Mart (DM) → Smaller version of DW



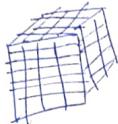
### Datawarehouse Operation

#### ETL (Extract Transform Load)



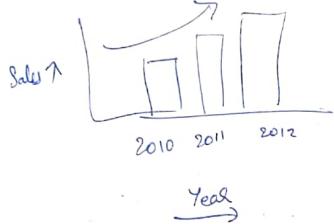
→ Extract from n sources. Transformation do work  
consistency among the data.

Transform: Slice, dice operations on data cube



Application of DW

→ Sales bed / bed in last 3 years



→ SAP software

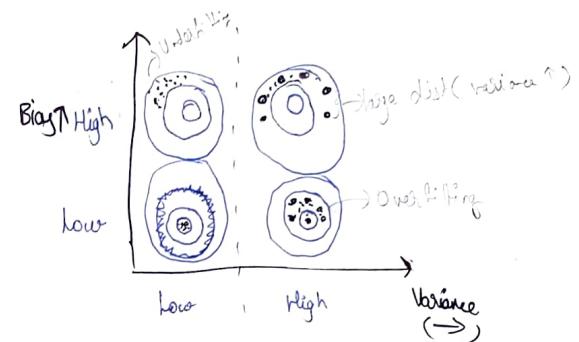
28/3

→ Difference b/w predicted value and actual value is bias error.

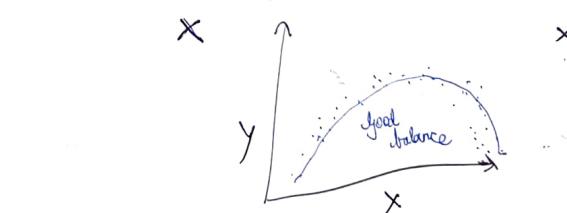
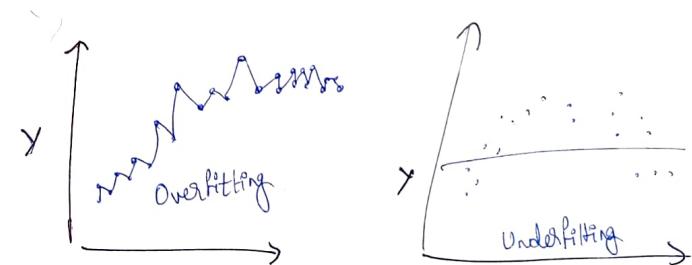
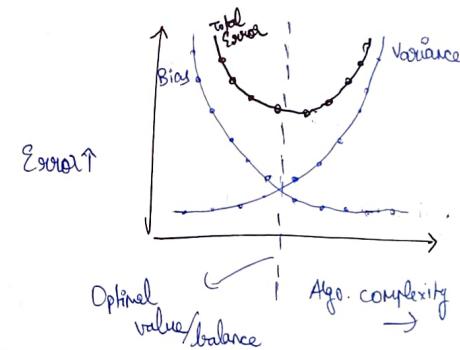
→ High bias, low variance → underfitting

→ Low bias, High variance → Overfitting

Bias / Variance trade off



$$\text{Total Error} = \text{Bias}^2 + \text{variance} + \text{irreducible error}$$



→ In supervised learning, underfitting happens when the model is unable to capture underlying pattern of data.

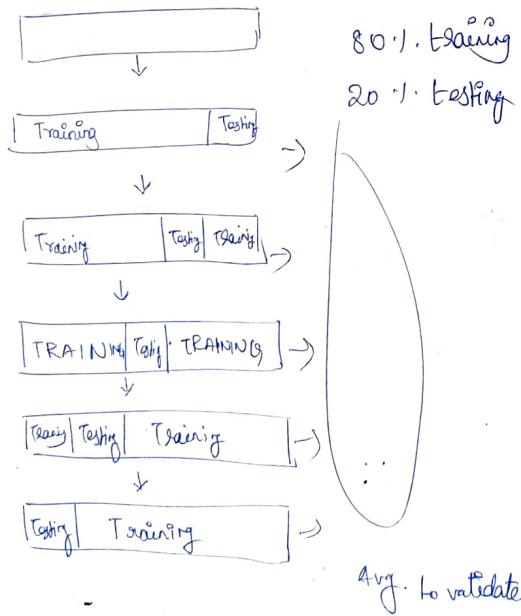
In creating a linear model using non-linear data,

underfitting can happen.

→ In supervised learning, overfitting happens when the model captures noise along with underlying pattern in data.  
↳ Happens when we train the model lot over a noisy dataset.

### k-fold cross validation

$$k = 1 \dots 5$$



→ LOOC  
→ Stratified

a. Total instances: 25  
k value: 5

→ each time 20/1 taken as test data.

$$20/1 \text{ of } 25 = \frac{20}{100} \times 25 = \frac{20}{5} = 5.$$

No. iteration	Training set observation	Testing set observation
1	All remaining values [5 6 7 ... 24]	[0 1 2 3 4]
2	[0 ... 4 10 ... 24]	[5 6 7 8 9]
3	[0 ... 9 15 ... 24]	[10 ... 24]
4	[0 ... 14 20 ... 24]	[5 ... 19]
5	[0 ... 19]	[20 ... 24]

White, the  
actual values  
not ...

You can divide available data into multiple folds or subsets. Using one of the folds or subsets as a validation set & training the models on remaining folds. This process repeated for multiple times. Each time using different fold as a validation set. Finally results from each validation steps are averaged to produce more robust estimate of the model's performance.

16 groups for iris dataset with k=5

	Sepal length	Sepal width	Petal length	Species
0	5.1	3.5	1.4	Virginica
1	4.9	3.0	1.4	Virginica
2	4.7	3.2	1.3	Vericolor
3	4.6	3.1	1.5	Virginica
4	4.5	2.3	1.0	Vericolor
5	4.5	2.3	1.5	Vericolor
6	4.5	2.0	1.6	Vericolor
7	4.5	1.8	1.7	Vericolor
8	4.3	1.6	1.4	Vericolor
9	4.3	1.5	1.0	Vericolor
10	5.8	2.7	5.1	Virginica
11	5.1	3.0	3.5	Virginica
12	5.9	3.0	4.2	Virginica
13	6.4	3.2	4.5	Virginica
14	6.5	3.2	4.5	Virginica
15	6.5	3.2	4.9	Virginica
16	6.5	3.2	4.9	Virginica

No.	Social no.	CGPA	Interacting	Practical knowledge	Communication skills	Job offer
0	1	>=9	Yes	Very good	Good	Yes
1	2	>=8	No	Good	Moderate	Yes
2	3	>=9	No	Average	Poor	No
3	4	<8	No	Average	Good	No
4	5	>=8	Yes	Good	Moderate	Yes
5	6	>=9 (6)	Yes	Good	Moderate	Yes
6	7	<8	Y	Good	Poor	N
7	8	>=9	N	Very good	Good	Y
8	9	>=8	Y	Good	Good	Y
9	10	>=8	Y	Avg.	Good	Y

$$20\% \text{ of } 10 = \frac{20}{100} \times 10 = 2$$

Training data

Testing data

Training  
data

Iteration

[0, 1]

{2, ..., 9}

2

[2, 3]

{0, 1, 4, ..., 9}

3

[4, 5]

{0, ..., 3, 6, ..., 9}

4

[6, 7]

{0, ..., 5, 8, 9}

5

[8, 9]

{0, ..., 7}

Apply k-fold

Iteration

Training  
data

Testing  
data

1

[0, 1]

[3, 4, 5, 6, 7, 8, 9]

2

[2, 3]

[0, 1, 4, 5, 6, 7, 8, 9]

3

[4, 5]

[0, 1, 2, 3, 6, 7, 8, 9]

4

[6, 7]

[0, 1, 2, 3, 4, 5, 8, 9]

5

[8, 9]

[0, 1, 2, 3, 4, 5, 6, 7]



## Noisy data

→ blank values

→ outlier data

→ inconsistent data (like ranking does not 'good', 'very good', 'most')

If value is not present,

→ can be removed.

→ fill mean of the existing values.

→ fill mean of the non-numeric values / mode of for non-numeric values

/ mode of for handle noise

→ Regression can be used to handle noise

→ Replace inconsistent values / replace with most frequent value

→ Replace inconsistent values / replace with most frequent value

Bin methods → 0 to n bins, calc avg for each bin

→ Centroid method to classify  
into classes

30/4

## Normalisation

1. Decimal scaling ( $\div$  by  $10^j$ ) [ $j=3$  if max val = 700]

2. Min-max normalisation

3. Z-score normalisation

## Decimal Scaling

$$V'_i = \frac{V_i}{10^j} \quad ; \text{ depends on max. value in dataset}$$

-10, 201, 301, -401, 501, 601, 701

Max abs. value = 701  $\Rightarrow j=3$

$$V'_i = \frac{-10}{10^3}, \frac{201}{10^3}, \frac{301}{10^3}, \frac{-401}{10^3}, \frac{501}{10^3}, \frac{601}{10^3}, \frac{701}{10^3}$$
$$= -0.01, 0.201, 0.301, -0.401, +0.501, 0.601, 0.701$$

## Min max normalisation

$$V' = \frac{V - \min(A)}{\max(A) - \min(A)} \quad (\text{new\_max}(A) - \text{new\_min}(A)) + \text{new\_min}(A)$$

## Z-Score normalisation

$$V' = \frac{V_i - \bar{A}}{\sigma_A}$$

For a normal dist,  $\bar{A}=0$ ,  $\sigma_A=1$

## Advantages of using normalisation

→ Improve performance of ML algo.

→ Helps to reduce outlier data / better handling of outlier data.

→ Improve interpretability of results.

→ Better generalisation.

## Disadvantages

→ Loss of information

→ Impact on outliers, interpretability

→ Additional computational cost

→ Pandas, numpy, sklearn can be used to normalize data.

## Normalization implementation

```

from sklearn import preprocessing
import pandas as pd
from sklearn.datasets import fetch_housing
housing = fetch_housing(as_frame=True)
scaler = preprocessing.MinMaxScaler()
d = scaler.fit_transform(housing.data)
scaled_df = pd.DataFrame(d, columns=housing.data.columns)
print(scaled_df)
    
```

to min max range: `min_maxScaler(feature_range=(0, 1))`

Median income	House age	Average no. of rooms per household	Avg bedrooms	Population	Avg occupants	Habitable area
1.0793	1.5606	0.087			3.788	122
1.0760	0.7843	0.076				
0.9320	2.00	0.105				
0.7093	2.0	0.070				
0.4695	2.0	0.077				

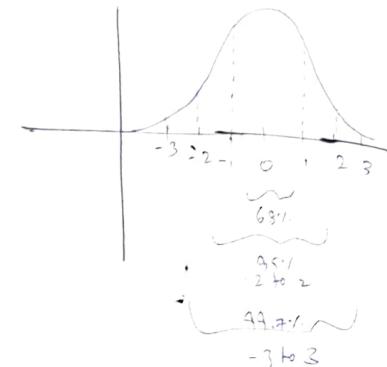
After min-max scaled ↓

1.13	0.423
1.1307	0.442
1.1255	0.428
1.1285	0.418
1.1285	0.418

6/4

Methods to detect outliers & removal

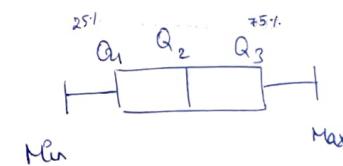
### 1) Standard deviation



→ Beyond -3 or +3 is outlier data.  
→ Within -1 to +1: Normal data

### 2) Interquartile range

#### 2) Box plot



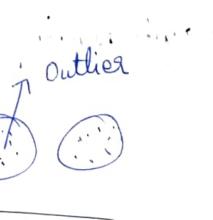
$$IQR = Q_3 - Q_1$$

Outliers  $\left\{ \begin{array}{l} < Q_1 - 1.5 * IQR \\ > Q_3 + 1.5 * IQR \end{array} \right.$

### 3) DBScan clustering

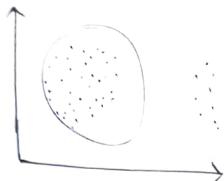
→ based on distance b/w centroids

→ very large dist: outlier



- 4) Isolation Forest
- 5) Random Cut Forest  
→ used by Amazon
- } specialized techniques

### Isolation forest



→ derived from original forest

### Random cut forest

→ based on low & high values.

↓  
Normal  
 $2000 - 26k$   
↓  
outliers  
 $210k$

e.g. Credit card purchases

→ Threshold like  $\geq 50,000$

sklearn, numpy can be used. (Anomaly)

### Applications

- 1) Quality assurance
- 2) Salary (not in expected range for given category)

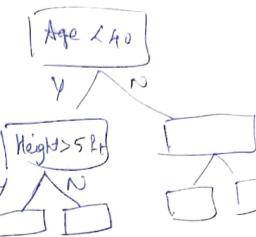
11/4

### Machine learning



- 1) Supervised (labelled data) - Decision tree, KNN, SVM, Random forest
- 2) Unsupervised (Unlabelled data)
- 3) Reinforcement (Agent - Reward)

### Decision tree



→ leaf nodes are class labels.

### K-nearest neighbour (KNN)

$$K = 1, \dots, n$$



Find dist. of test data pt. with all other data points. Find the closest  $K$  points & avg. of their weight values will be weight of test data pt.

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

2) Manhattan distance  $(\sum_{i=1}^k |x_i - y_i|)$  for continuous variables

3) Hamming distance  $(D_H) \sum_{i=1}^k |x_i - y_i|$  for categorical variables

$$x = y \Rightarrow D = 0$$

$$x \neq y \Rightarrow D = 1$$

Q.

ID	<u>Height</u>	<u>Age</u>	<u>Weight</u>
1	5	45	77
2	5.11	26	44
3	5.6	30	55
4	5.9	34	59
5	4.8	40	72
6	5.8	36	60
7	5.3	19	40
8	5.8	28	60
9	5.85	23	45
10	5.6	32	58

Testing: 11 5.5 38

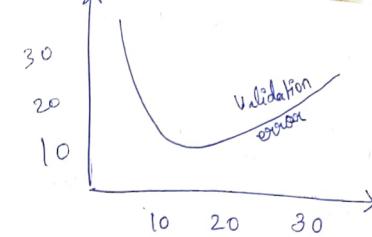
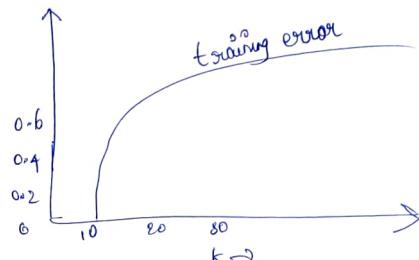
Distances b/w  
test & data points

- 1) 7.017
- 2) 12.007
- 3) 8.001
- 4) 4.02
- 5) 2.119
- 6) 2.023
- 7) 19.002
- 8) 10.005
- 9) 15
- 10) 6.001

Min dist points are 6, 5, 4

$$\text{Avg of the weights of these Pts.} = \frac{191}{3} = 63.67$$

ID	<u>Ht.</u>	<u>Age</u>	<u>Weight</u>
4	5.9	34	59
5	4.8	40	72
6	5.8	36	60



Name	<u>Age</u>	<u>Gender</u>	<u>Class of sport</u>	Dist from test
Ajay	32	0	Football	17.03
Mack	40	0	Neither	25.02
Sarah	16	1	Cricket	1
Zaira	34	1	Cricket	19
Sachin	55	0	Neither	40.013
Rahul	40	0	Cricket	25.02
Poopa	20	1	Neither	5
Smith	15	0	Cricket	1
Lakshmi	55	1	Football	40
Michael	15	0	Football	1

Test

Angelia 15 1 ?

Min dist Pts are Sarah, Smith, Michael

Cricket, Cricket, Football

$\Rightarrow$  Cricket.

## Steps used for supervised learning algo.

- 1) Determine type of training dataset
- 2) Collect / gather labelled training data
- 3) Split training dataset into training, test & validation dataset.  
↳ subset of training dataset
- 4) Determine input features of training dataset which should have enough knowledge which the model can accurately predict output
- 5) Determine suitable algorithm for the model such as SVM, decision tree, Random forest etc.
- 6) Execute algorithm on training dataset. Something we need validation set as a control parameter which are subset of training dataset.
- 7) Evaluate accuracy of your model by providing test dataset. If the model predict correct output, it means model is accurate.

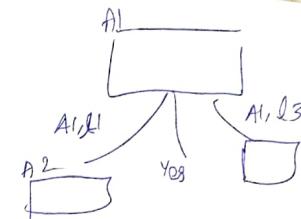
13/4

## Decision Tree

$$\text{Entropy, } E(S) = - \sum_{i=1}^n P(x_i) \log_2 P(x_i)$$

$$\text{Information gain, } IG(S, A) = E(S) - E(S, A)$$

Total entropy - Average entropy



$\Rightarrow 0 \rightarrow \text{No split}$   
 $\neq 0 \rightarrow \text{Further split.}$

Outlook	Temperature	Humidity	Wind	Play - Football
Sunny	Hot	High	Weak	N
Sunny	Hot	High	Strong	N
Overcast	Hot	High	Weak	Y
Rain	Mild	High	Weak	Y
Rain	Cool	Normal	Weak	Y
Rain	Cool	Normal	Strong	Y
Overcast	Cool	Normal	Strong	N
Sunny	Mild	High	Weak	Y
Sunny	Cool	Normal	Weak	N
Rain	Mild	Normal	Weak	Y
Sunny	Mild	Normal	Strong	Y
Overcast	Mild	High	Strong	Y
Overcast	Hot	Normal	Strong	Y
Rain	Mild	High	Weak	N

$$E(S) = - \left[ \frac{9}{14} \log_2 \left( \frac{9}{14} \right) + \frac{5}{14} \log_2 \left( \frac{5}{14} \right) \right]$$

$$= 0.94$$

Outlook

Outlook

Play Football		Total
Sunny	Yes	5
Overcast	No	4
Rainy	Yes	5
		14

$$E(S, \text{outlook}) = \frac{5}{14} \cdot E(3, 2) + \frac{4}{14} \cdot E(4, 0)$$

$$+ \frac{5}{14} \cdot E(2, 3)$$

$$= \frac{5}{14} \cdot -\left(\frac{3}{5} \cdot \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5}\right)$$

$$+ \frac{4}{14} \cdot 0 + \frac{5}{14} \left(-\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5}\right)$$

$$= 0.693$$

$$IG(S, \text{outlook})$$

$$\cancel{0.94 - 0.693} = 0.247$$

Temperature

## Play Football

Play Football		Total
Hot	Yes	2
Mild	No	2
Cool	Yes	1
		4
		14

$$E(S, \text{temp}) = \frac{4}{14} \cdot E(2, 2) + \frac{6}{14} \cdot E(4, 2) + \frac{4}{14} \cdot E(3, 1)$$

$$= \frac{4}{14} - \left(\frac{2}{4} \log \frac{2}{4} + \frac{2}{4} \log \frac{2}{4}\right) + \frac{6}{14} - \left(\frac{4}{6} \log \frac{4}{6} + \frac{2}{6} \log \frac{2}{6}\right)$$

$$+ \frac{4}{14} - \left(\frac{3}{9} \log \frac{3}{9} + \frac{1}{9} \log \frac{1}{9}\right)$$

$$= 0.911$$

$$IG(S, \text{temp}) = 0.94 - 0.911 = 0.029$$

(S, Humidity)

	Yes	No	Total
High	3	4	7
Normal	6	1	7

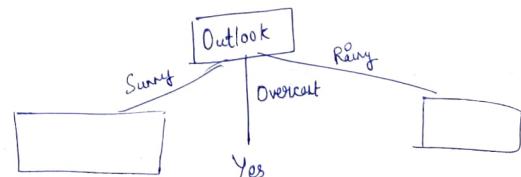
$$E(S, \text{humidity}) = \frac{7}{14} \cdot E(3, 4) + \frac{7}{14} \cdot E(6, 1) = 0.788$$

$$IG(S, \text{hum}) = \frac{7}{14} = 0.94 - 0.788 = 0.152$$

$$E(S, \text{wind}) = 0.8932$$

$$IG(S, \text{hum}) = 0.048$$

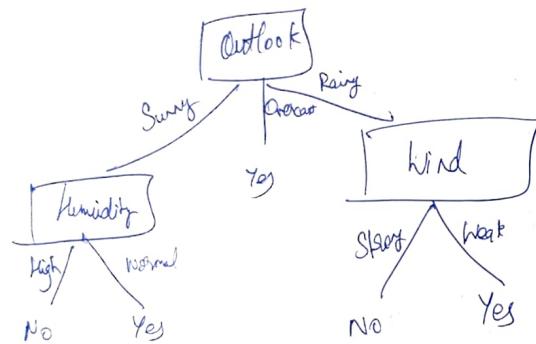
Highest IG is for humidity outlook



<u>Outlook</u>	<u>Temp</u>	<u>Humidity</u>	<u>Wind</u>	<u>Play Football</u>
Sunny	High	H	S	N
Sunny	H	H	W	N
Sunny	Mid	H	W	Y
Sunny	Cool	N	S	Y
Sunny	Mild	N		

$$E(S) = - \left[ \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right]$$

∴ Humidity has highest IG



$$\text{Gini}(E) = 1 - \sum_{j=1}^k P_j^2 \quad \text{Used for CART, Classification & Regression Tree}$$

18/4

For Outlook = Sunny

$$E(S) = - \left[ \frac{3}{5} \log \frac{3}{5} + \frac{2}{5} \log \frac{2}{5} \right] = 0.97$$

$$IG(S | \text{Outlook} = \text{Sunny}) = 0.97 - 1.732$$

Temp

<u>Temp</u>	<u>Yes</u>	<u>No</u>	Total
Hot	0	2	2
Cool	1	0	1
Mild	1	1	<u>2</u> 5

$$E(S, \text{Temp}) = - \frac{2}{5} \left( \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right)$$

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

$$IG(S, \text{Temp}) = 0.97 - 2/5 = 0.57$$

<u>Humidity</u>	<u>Yes</u>	<u>No</u>	Total
High	0	3	3
Normal	2	0	<u>2</u> 5

$$E(S, \text{Humidity}) = 0$$

$$\log_{10} a = \frac{\log_{10} a}{\log_{10} 10}$$

$$IG(S, \text{Humidity}) = 1.732$$

Gini

$$Gini(S) = 1 - \sum_{i=1}^2 P_i^2$$

$$= 1 - \left( \left(\frac{9}{14}\right)^2 + \left(\frac{5}{14}\right)^2 \right)$$

$$= 0.4591$$

$$Gini(S, Outlook) = \frac{5}{14} Gini(S_1, 2) + \frac{4}{14} Gini(S_1, 10) + \frac{5}{14} Gini(S_2, 3)$$

$$= \frac{5}{4} \left( 1 - \left( \left(\frac{8}{5}\right)^2 + \left(\frac{2}{5}\right)^2 \right) \right) + \frac{4}{14} \times 0 + \frac{5}{14} \times 1 - \left( \left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right)$$

$$= 0.171 + 0 + 0.171 = 0.342$$

$$Gini(S, Outlook) = 0.459 - 0.342 = 0.117$$

$$Gini(S, Temperature) = 0.459 - 0.4405 = 0.0185$$

$$Gini(S, Humidity) = 0.459 - 0.3674 = 0.0916$$

$$Gini(S, Windy) = 0.459 - 0.4286 = 0.0304$$

Apply decision tree for (GPA, using Information gain  
 Job after dataset)  
 ↪ back page

$$E(S) = - \left[ \frac{4}{10} \log \frac{4}{10} + \frac{7}{10} \log \frac{7}{10} \right] = 0.266 \rightarrow 0.882$$

<u>GPA</u>	<u>Yes</u>	<u>No</u>	<u>Total</u>
$\geq 9$	3	1	4
$\geq 8$	4	0	4
$< 8$	0	2	<u><math>\frac{2}{10}</math></u>

$$E(S, GPA) = \frac{4}{10} E(S_1, 1) + \frac{4}{10} \times 0 + \frac{2}{10} \times 0 = \frac{4}{10} \left( \frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \right) = 0.3246$$

$$IG(S, GPA) = 0.5574$$

### Interactive

		<u>Yes</u>	<u>No</u>	<u>Total</u>
<u>Yes</u>	<u>Yes</u>	5	1	6
	<u>No</u>	2	2	4
		<u>10</u>		

$$E(S, \text{Interactive}) = -\frac{6}{10} \left( \frac{5}{6} \log \frac{5}{6} + \frac{1}{6} \log \frac{1}{6} \right)$$

$$+ -\frac{4}{10} \left( 2 \cdot \frac{2}{4} \log \frac{2}{4} \right)$$

$$= -\frac{6}{5} (0.65) + \frac{4}{10}$$

$$= 0.39$$

$$IG(S, \text{Interactive}) = 0.492$$

(GPA is High)

GPA  
Yes      No  
Yes      Yes

### Practical knowledge

		<u>Yes</u>	<u>No</u>	<u>Total</u>
<u>Very g</u>	<u>Yes</u>	2	0	2
	<u>No</u>	4	1	5
	<u>avg</u>	1	2	3
		<u>10</u>		

$$E(S, \text{Practical know}) = \frac{5}{10} E(1,1) + \frac{3}{10} E(1,2)$$

$$= -0.5 \left( \frac{4}{5} \log \frac{4}{5} + \frac{1}{5} \log \frac{1}{5} \right) - 0.3 \left( \frac{1}{3} \log \frac{1}{3} + \frac{2}{3} \log \frac{2}{3} \right)$$

$$= 0.361 + 0.276 = 0.637$$

$$IG = 0.2456$$

### Comm. skills

		<u>Yes</u>	<u>No</u>	<u>Total</u>
<u>Good</u>	<u>Yes</u>	4	1	5
	<u>Moderate</u>	3	0	3
	<u>Poor</u>	0	2	2
		<u>10</u>		

$$E(S, \text{Comm}) = -\frac{5}{10} \left( \frac{4}{5} \log \frac{4}{5} + \frac{1}{5} \log \frac{1}{5} \right)$$

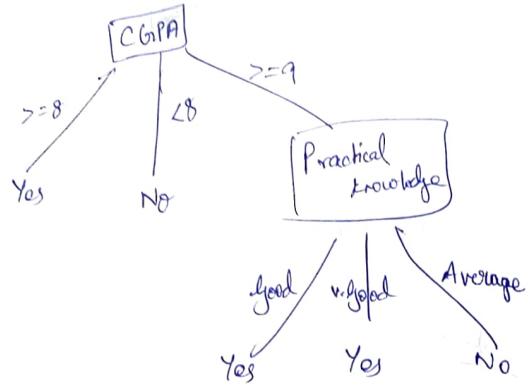
$$= 0.361$$

$$IG = 0.521$$

### GPA ≥ 9

<u>GPA</u>	<u>Interactive</u>	<u>Practical know</u>	<u>Comm skills</u>	<u>Total score</u>
<u>&gt;= 9</u>	<u>Y</u>	<u>VG</u>	<u>G</u>	<u>Y</u>
<u>&gt;= 9</u>	<u>N</u>	<u>A</u>	<u>P</u>	<u>N</u>
<u>&gt;= 9</u>	<u>Y</u>	<u>G</u>	<u>M</u>	<u>Y</u>
<u>&gt;= 9</u>	<u>N</u>	<u>VG</u>	<u>G</u>	<u>Y</u>

$$E(S) = - \left( \frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \right) = 0.8113$$



<u>Intrative</u>	Yes	No	
Y	2	0	2
N	1	1	<u>2</u>
			<u>4</u>

$$-\frac{2}{4} \left( \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) = \frac{1}{2} \log_2 \frac{1}{2}$$

$$= -\frac{1}{2} \log_2 2^1 = 1/2 = 0.5$$

$$IG_I = 0.3113.$$

<u>Prac</u>	Yes	No	Total
VG	2	0	2
G	1	0	1
A	0	1	<u>1</u>

$$E = 0, IG_I = 0.8113.$$

<u>Comm</u>	Yes	No	Total
G	2	0	2
M	1	0	1
P	0	1	<u>1</u>

Prac know = VG.

$$\begin{array}{c} \text{Prac know} = VG \\ \hline \text{Pr } K \\ \hline \frac{VG}{VG} \\ \hline \end{array} \quad \begin{array}{c} \text{Job} \\ \hline Y_N \\ \hline Y_S \\ \hline \end{array}$$

Prac k = G

$$\begin{array}{c} \text{Prac k} = G \\ \hline \text{Pr } K \\ \hline \frac{G}{G} \\ \hline \end{array} \quad \begin{array}{c} \text{Job} \\ \hline Y \\ \hline \end{array}$$

Prac k = Avg

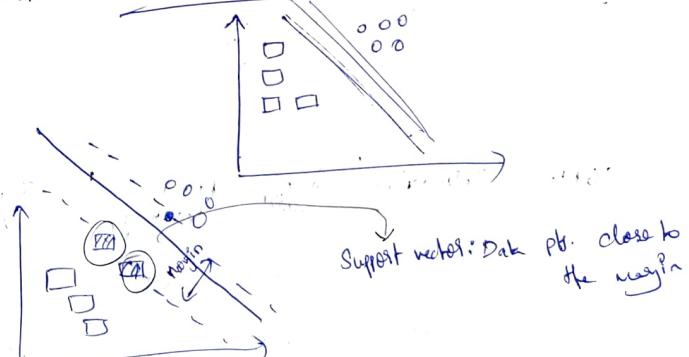
$$\begin{array}{c} \text{Prac k} = \text{Avg} \\ \hline \text{Pr } K \\ \hline \text{Avg} \\ \hline \end{array} \quad \begin{array}{c} \text{Job} \\ \hline N_O \\ \hline \end{array}$$

CGPA < 8

$$\begin{array}{c} \text{CGPA} < 8 \\ \hline \text{Job} \\ \hline \text{No} \\ \hline \text{No} \\ \hline \end{array}$$

$$\begin{array}{c} \text{CGPA} > 8 \\ \hline \text{Job} \\ \hline Y \\ \hline Y \\ \hline Y \\ \hline Y \\ \hline \end{array}$$

SUPPORT VECTOR MACHINE



Support vector: Data pt. close to the margin

Cost Function

$$C(x, y, f(x)) = \begin{cases} 0 & \text{if } y \cdot f(x) \geq 1 \\ 1 - y \cdot f(x) & \text{else} \end{cases}$$

gradient update

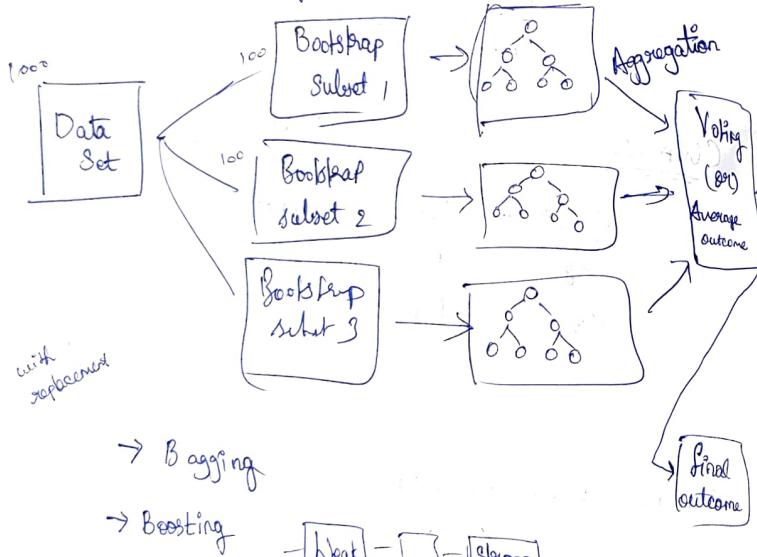
$$w = w - \alpha (2xw) \rightarrow \text{No misclassification}$$

$$w = w + \alpha (y_i * x_i - 2xw) \rightarrow \text{final gradient update for proper classification}$$

20/4

## Random Forest

→ Ensemble algorithm



→ Bagging

→ Boosting



Weak to Strong classifier

Q.

CGPA	Interactions	Comm. skills	Prac. knowledge	Job offer
1. $\geq 9$	Y	G <sub>1</sub> and G <sub>2</sub>	G <sub>1</sub>	Y
2. $< 9$	N	Moderate	G <sub>2</sub>	Y

	$\geq 9$	No	M	Avg	N
1.	$< 9$	N	M	A	N
5.	$\geq 9$	Y	M	G	Y

Bootstrap Subsets S1

1.	$\geq 9$	Y	G	f	Y
2.	$< 9$	N	M	G	Y
3.	$\geq 9$	N	M	A	N
4.	$\geq 9$	N	M	A	N
5.	$\geq 9$	Y	M	G	Y

S2

1.	$< 9$	N	M	G	Y
2.	$\geq 9$	N	M	A	N
3.	$\geq 9$	N	M	A	N
4.	$\geq 9$	Y	M	G	Y
5.	$\geq 9$	Y	M	G	Y

S3

1.	$\geq 9$	Y	G	G	Y
2.	$\geq 9$	Y	G	G	Y
3.	$< 9$	N	M	G	Y
4.	$\geq 9$	N	M	A	N
5.	$\geq 9$	N	M	A	N

Chapter 2) Constant decision rule for cash inflow.

3) Choose no. of features to be used in each tree. (2)

4) choose depth of each decision tree to see 2

5) P + H P, Many steps for 3 decision tree

Repete the following words.  
What are conjectures.

a) Choose book sample of 5 data items  
Decision tree 1.

a) How to construct decision rule

b) Compute the best split feature to

b) compare it, be placed at each node.

c) Construct decision tree using CART

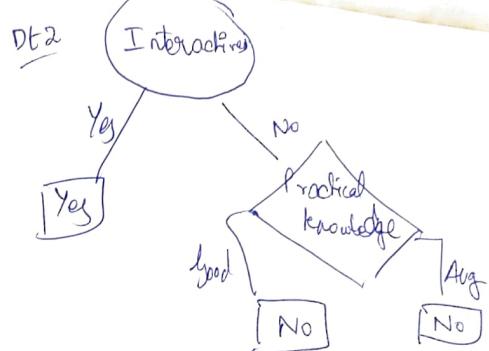
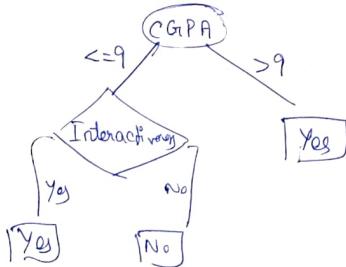
algorithm

6. Combine predictions from all the classifiers

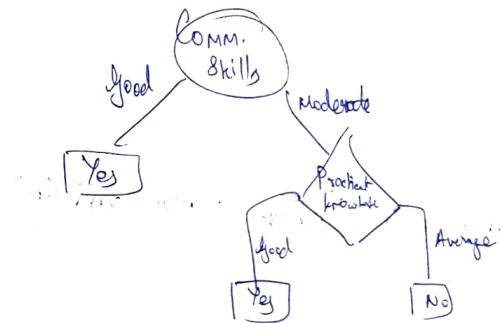
v. Combine or ensemble -  
taking an average for regression problem or majority voting for classification.

1st ~~test~~ using CGPA & interactiveness.

dt 1



DS Comm & Pk



4<sup>th</sup> instance as test date

DT1: Yes, DT2: No, DT3: No

Afghan Voting: No

## Advantages:

- Can solve both classification & regression
- Reduces overfitting & variance compared to single decision tree.
- Provides more accuracy
- Performs well even with small data set

Digitized by srujanika@gmail.com

Dig: → Inherent complexity in construction & prediction  
→ more computational resources

27/4

Types

1. Partition based clustering / Centroid based clustering

2. Hierarchical based / Connectivity based clustering

3. Density based / Model based

4. Distribution based clustering

5. Fuzzy based clustering

Partition based Clustering (k-means clustering)

Object	x-coordinate	y-coordinate
1	2	4
2	4	6
3	6	8
4	10	4
5	12	4

$\underline{c_1} \quad \underline{c_2}$   
 $(4, 6) \quad (12, 4)$

Iteration 1

$$\text{Dist}(1, \text{centroid } 1) = 2\sqrt{2}$$

$$\text{Dist}(1, \text{centroid } 2) = 10$$

$$\text{Dist}(3, \text{centroid } 1) = 2\sqrt{2}$$

$$\text{Dist}(3, \text{centroid } 2) = \sqrt{52}$$

$$\text{Dist}(4, c_1) = \sqrt{40}$$

$$\text{Dist}(4, c_2) = \sqrt{48}$$

After, Iteration 1

c1

(4, 6)

(2, 4)

(6, 8)

c2

(10, 4)

(12, 4)

(11, 4)

New centroid : (4, 6)

Iteration -2

$$\text{Dist}(2, c_1) = 0$$

$$\text{Dist}(2, c_2) = \sqrt{53}$$

$$\text{Dist}(3, c_1) = \sqrt{8}$$

$$\text{Dist}(3, c_2) = \sqrt{81}$$

$$\text{Dist}(4, c_1) = \sqrt{40}$$

$$\text{Dist}(4, c_2) = \sqrt{48}$$

$$\text{Dist}(5, c_1) = \sqrt{68}$$

$$\text{Dist}(5, c_2) = \sqrt{41}$$

$$\text{Dist}(6, c_1) = \sqrt{5}$$

$$\text{Dist}(6, c_2) = \sqrt{51}$$

After Iteration 2

c1

(2, 4)

(4, 6)

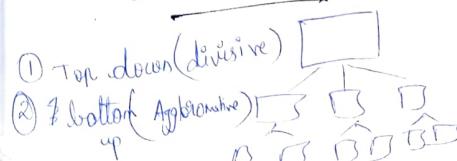
(6, 8)

c2

(10, 4)

(12, 4)

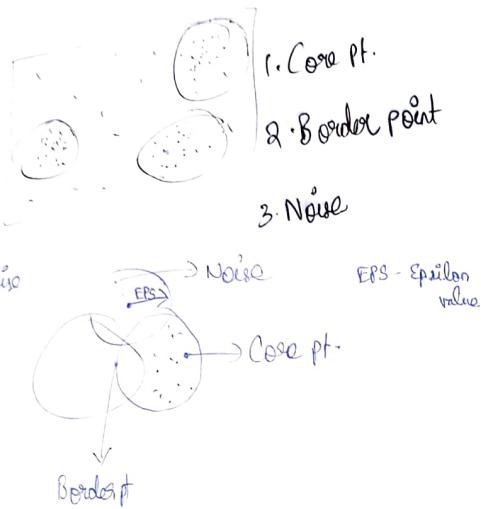
Hierarchical based



→ Clusters & sub clusters

Specify based

DB Scan  
DB SCAN  
Density Based  
Spatial Clustering  
Application with Noise



## BIRCH

Balanced Iterative Reducing and Clustering using Hierarchies

BR

1. CF Tree (Clustering Feature tree)
2. Threshold

Q5 Use k-means algo ( $k=2$ ) and show the result for the following

S.no	X	Y
1	3	5
2	7	8
3	12	5

$$\begin{array}{c|c|c} 4 & 16 & 9 \end{array}$$

Let  $C_1 = (3, 5)$ ,  $C_2 = (12, 5)$

$$\begin{array}{c|c} C_1 & C_2 \\ \hline (3, 5) & (12, 5) \\ (7, 8) & (16, 9) \end{array}$$

$$\text{Dist}(2, C_1) = \sqrt{25}$$

$$\text{Dist}(2, C_2) = \sqrt{84}$$

$$\text{New } C_1 = (5, 6.5)$$

$$\text{New } C_2 = (14, 7)$$

$$\begin{array}{c|c} C_1 & C_2 \\ \hline (3, 5) & (12, 5) \\ (7, 8) & (16, 9) \end{array}$$

$$\text{Dist}(2, C_1) = 2.5$$

$$\text{Dist}(2, C_2) = 7.07$$

## Hadoop

→ Distributed File System (HDFS)

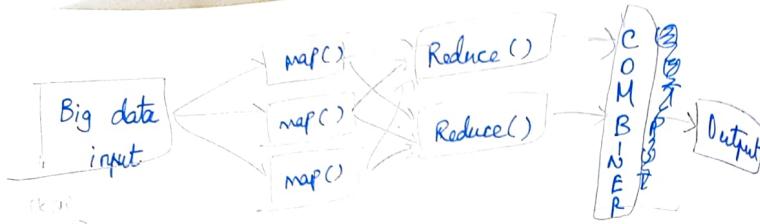
→ Distributed Processing (MapReduce)

→ Size of HDFS block 128-256 MB

→ Default replication factor - 3

Important components:

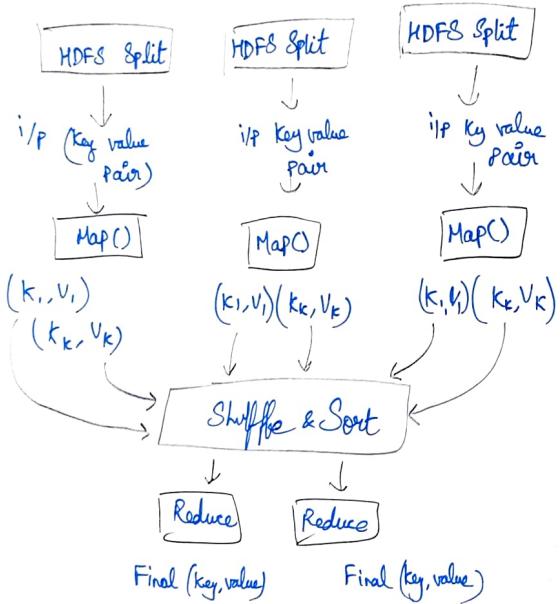
- 1) HDFS
- 2) Map Reduce
- 3) YARN framework
- 4) Hadoop common utilities (pig, hbase, hive)



Reduce operations:

- 1) Sort()
- 2) Shuffle()

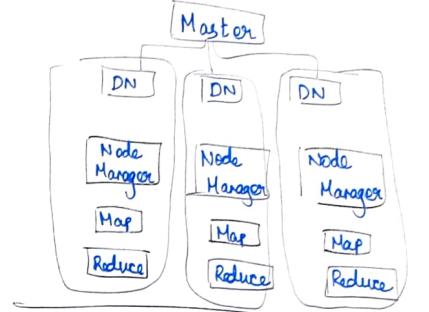
$k_i = \text{Key}_i$     $kk = \text{key}_k$     $V_k = \text{Value}_k$



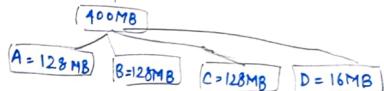
### HDFS

- Based on Google FS (GFS)
- Size of each block 128 - 256 MB
- Name node & data node
  - (Master) (Stores file blocks)
  - (Stores metadata)

→ Rack aware e.g. Rack 1 has data nodes DN1 - 10  
Rack 2 " DN11 - 20



### Data blocks in HDFS



$$\text{Rep. factor} = 3 \Rightarrow \text{Total file blocks} = 3 * 4 = 12 \text{ blocks}$$

→ Usually 500 DN are used.

- Advantages
- 1) Size of unix file blk is 4 kB while HDFS has 128 MB
  - 2) Replication
  - 3) High throughput
  - 4) Commodity hardware (works with inexpensive hardware)
  - 5) High Fault tolerant
  - 6) Software managing failures

### Installation

1. Pseudo distributed
2. Fully distributed (Min 2 DN present)

### XML files

Core-Site.xml  
Yarn-Site.xml  
HDFS-Site.xml  
Map Red-Site.xml

} Config files

pig  
hive  
hbase  
sqoop

Data Science  
 Analysis

9/5

## Hive - External Table

```
> CREATE EXTERNAL TABLE emp.employee_external
  ( id int, name String, age int, gender String)
ROW FORMAT DELIMITED FIELDS TERMINATED BY
';' LOCATION '/user/hive/data/employee_external'
```

## Temporary Table

```
> CREATE temporary TABLE emp.employee_temp(
  id int, ...) ROW ... BY '...';
```

## Partitioning Table

```
> CREATE TABLE zipcodes (RecordNumber int, Country
String, city String, zipcode int) PARTITIONED BY
(State String) ROW FORMAT ... ;
```

> DESCRIBE zipcodes;

> DROP TABLE <tablename>

> DELETE FROM <tablename> WHERE ...

## MongoDB

→ Document based, dynamic schema

C++

Cross platform

High performance / Scalable / Availability

RDBMS	MongoDB
DB	Column, Field
Table	DB
Tuple	Embeddable doc

RDBMS	MongoDB
Column, Field	Primary, -id
Join	Collection
Embedded doc	Document

→ index on any attribute ; replication, auto sharding, such queries

### Data Types

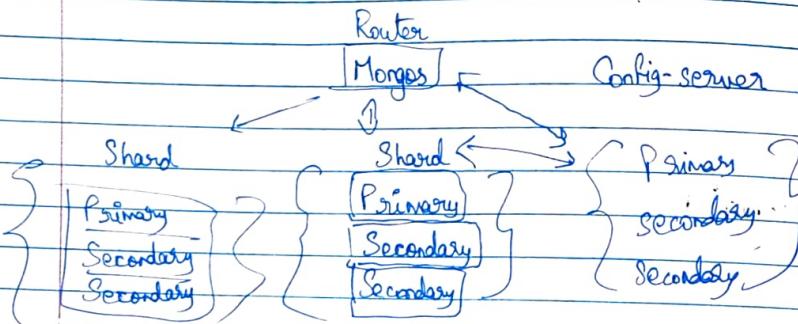
String, Integer, Boolean, double, min/max keys, arrays, timestamp, object, null, symbol, date, object ID, binary data, code, greg ex

db. find()

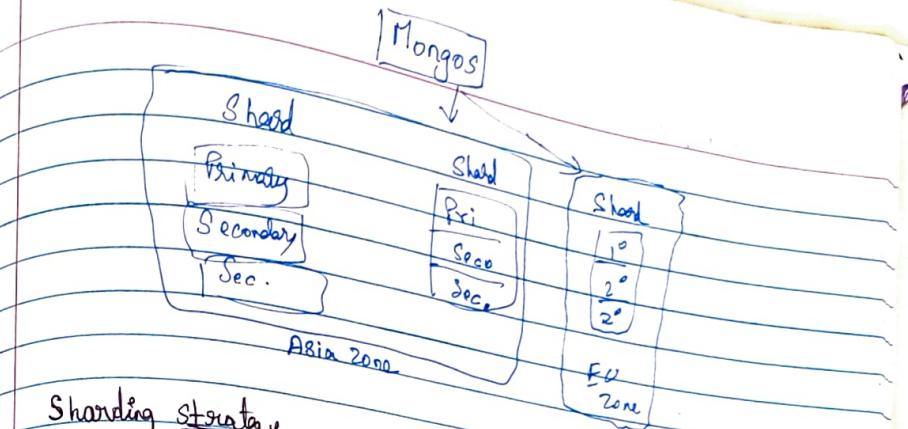
.find() .pretty()

## 115 Positioning in MongoDB

### Sharding



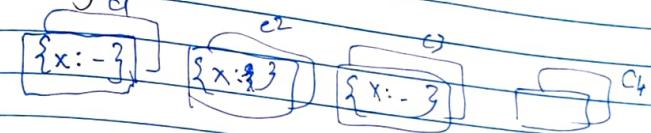
### Zone Sharding



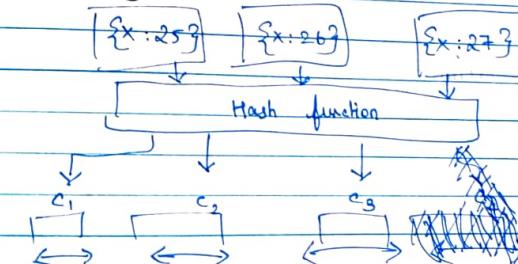
### Sharding Strategy

- Ranged Sharding
- Hashed Sharding

#### Ranged Sharding



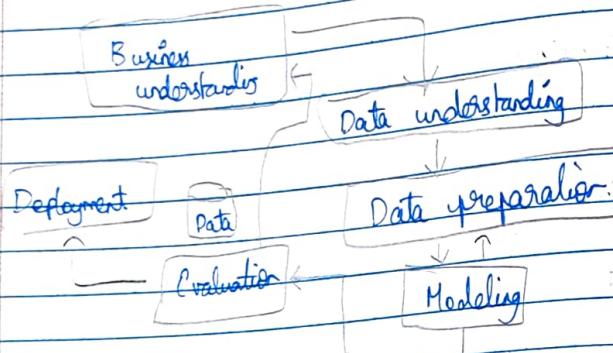
#### Hashed Sharding



## Benefits

- Increased Read/Write throughput
- Increased storage capacity
- Data locality

## 13/5 Data Munging (Wrangling)



1. Pre processing
2. Standardize dataset into understandable format
3. Cleaning data from noise & erroneous elements
4. Consolidate data from various sources.
5. Matching data with existing dataset.
6. Filtering data through determined settings.

## Data cleaning

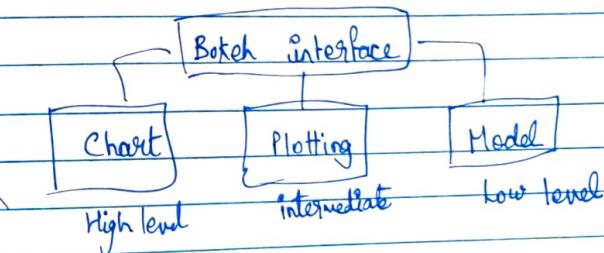
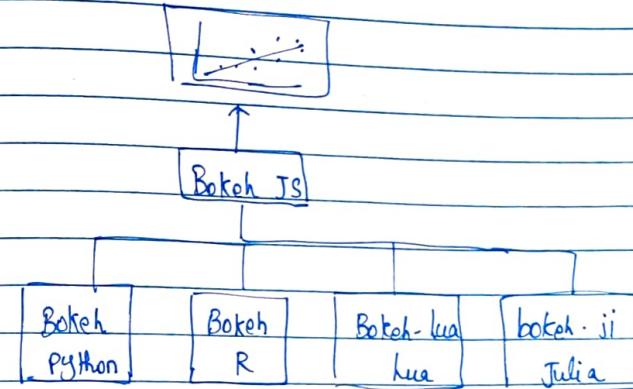
- Data audit
- Workflow specification & execution
- Post processing control

## Tools for data munging

- Excel Power Query / spreadsheets
- OpenRefine
- Google Data Prep
- Talend
- Data wrangler
- CSV kit

## 16/5 Data Pipeline

Interactive data visualization using Bokeh



## Methodology to create chart

1. Import the library and functional methods
2. Prepare a data
3. Set output mode (notebook / web browser / server)
4. Create chart with styling options (server)
5. Visualize the chart

## Bar chart on web browser

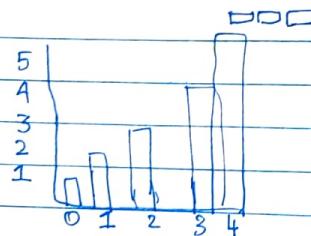
```
from bokeh.charts import Bar, output_file, show
data = {'y': [1, 2, 3, 4, 5]}
```

```
output_file('bar.html', title='BarChart Example')
```

```
P = Bar(data, title='Bar chart example',
        xlabel='x', ylabel='values', width=400,
        height=400)
```

Show (P)

## Output

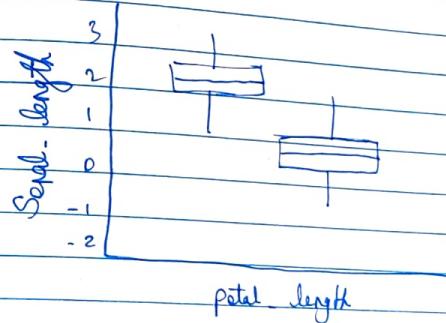


Compare distribution of sepal length & petal length of iris dataset using box plot on notebook screen.

```
from datasets import load_iris
import pandas as pd
df = pd.DataFrame(iris)
```

```
df.columns = ['Petal width', 'Petal length', 'Sepal width',
              'Sepal length']
from bokeh.charts import Boxplot, output_notebook
data = df[['Petal length', 'Sepal length']]
output_notebook()
P = Boxplot(data, width=400, height=400)
show(P)
```

## Output

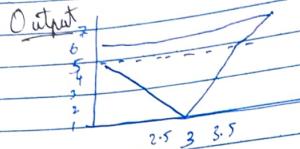


Create line plot on bokeh server

> bokeh-server

```
from bokeh.plotting import figure
out_server('line') output_server, show
```

`p = figure (Plot_width=400, Plot_height=400)`  
`p.line ([5,2,3,4,5], [5,7,2,4,5], line_width=2)`  
`show(p)`



Create a square mask on xy frame of notebook  
Combining line plots

from batch\_plotting import figure, output\_notebook,  
show

`p = figure (Plot_width=400, Plot_height=400)`

`p.square ([2,5,6,4], [2,3,2,1,2],`

`size=20, color='navy')`

`show(p)`

Output

