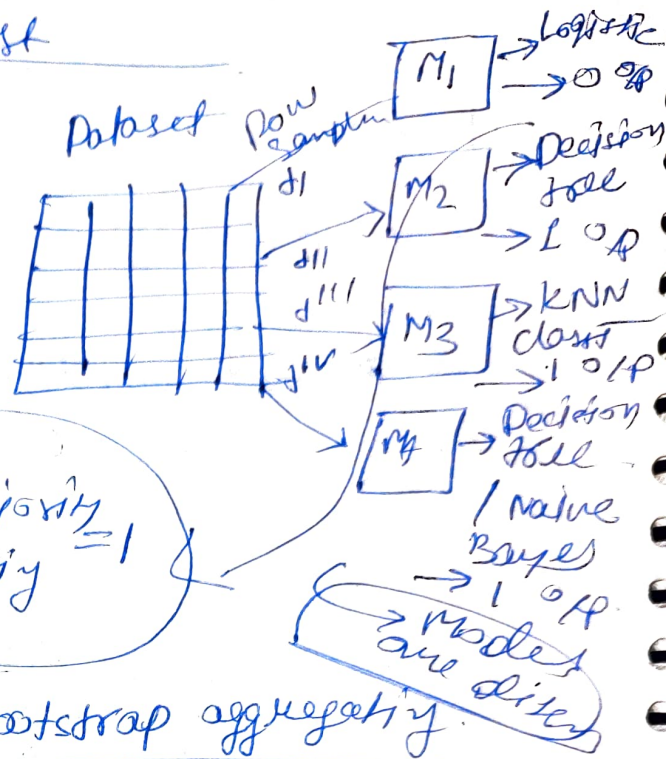
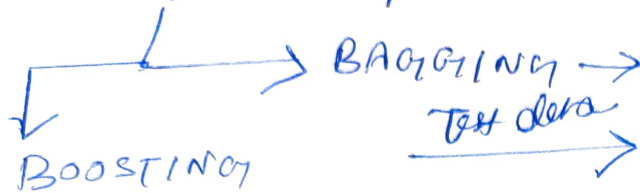


Adaboost & Random Forest

1

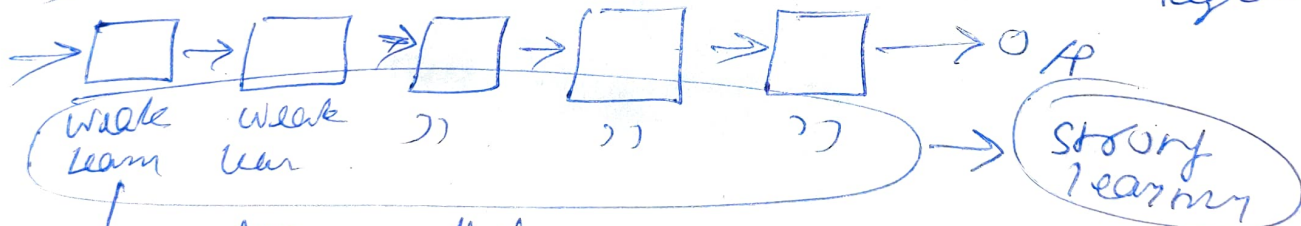
Ensemble Techniques:-



Boosting

2

Sequential models



Mean prediction is very weak that's why we combined models.

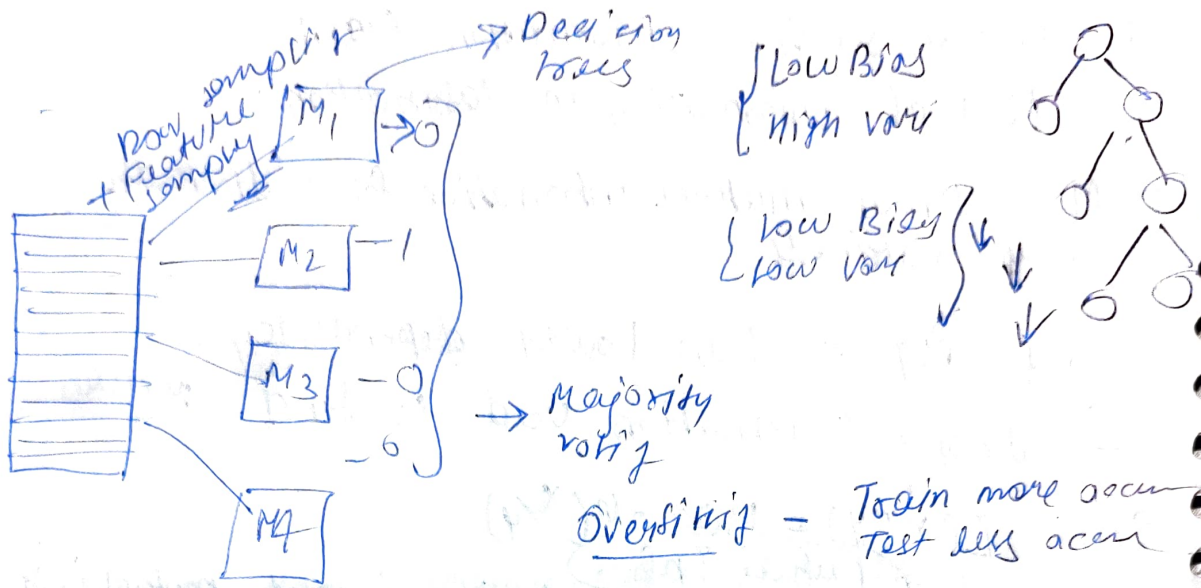
In Regression O/P mean will be taken

- Boosting
- 1 Random F classification
 - 2 Random F Regression ✓

- 1 Adaboost
- 2 Gradient
- 3 Xgboost

Topics Data Scientist

DT



Is normalization required in random forest or DT → No, because we are using split so some or norm do not that much necessary.

Is random forest impacted by outliers → Yes coz they hinders splitting process.

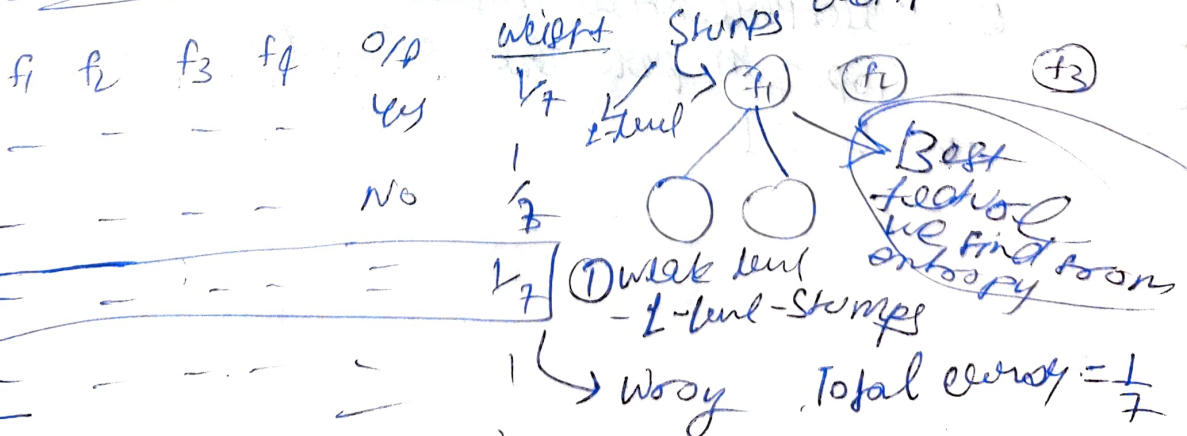
Is KNN impacted by outliers → Yes.

Boosting Technique → Some in sequential

③ Adaboost

Overall = 1

Information gain



② Performance of stumps:

$$= \frac{1}{2} \log_2 \left(\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right) = \frac{1}{2} \log_2 \left(\frac{1 - \frac{1}{7}}{\frac{1}{7}} \right) = 0.895$$

Deep learning & NLP

③ New sample weight = weight $\times e^{-P_s^{(2)}}$ $= \frac{1}{7} \times e^{-0.895} = 0.09$

correct

Incorrect record = weight $\times e^{P_s} = \frac{1}{7} \times e^{0.895} = 0.349$

New weight

Is it $\Sigma = 1$.

Normalised weight

Bucket

0.05 \div 0.649

0.05

0.05

0.05

0.05

0.05 > 0.349

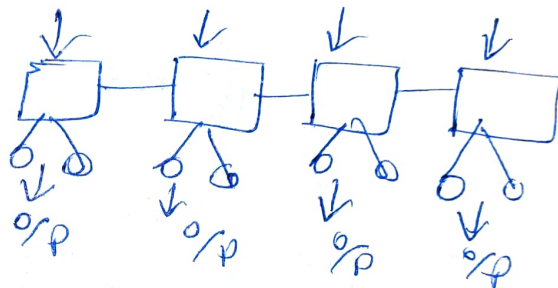
0.05

0.05

0.05

0.649

0.07	[0 - 0.07]
0.07	[0.07 - 0.14]
0.07	[0.14 - 0.21]
0.537	[0.21 - 0.75]
0.07	[0.75 - 0.82]
0.07	
0.07	
0.07	



Regression = Avg
classification = Majority

④ Boost

Black box vs White Box Model

Linear Regression \rightarrow White \rightarrow We can visualize

Random Forest \rightarrow Black

Decision Tree \rightarrow White

RNN \rightarrow Black box.

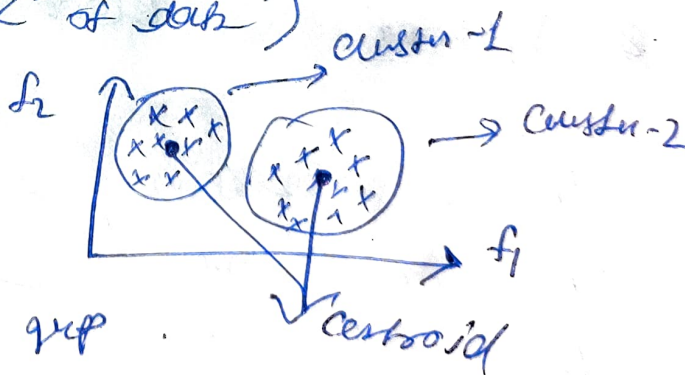
Day-6

Unsupervised ML

Don't have specific %
Clusters

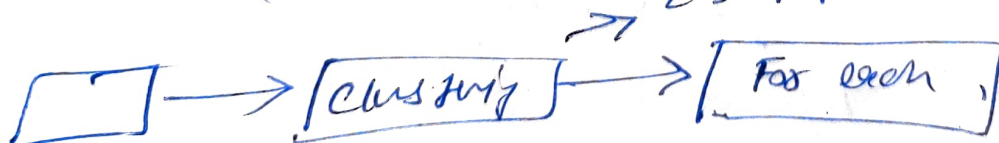
- ① k-means clustering
- ② Hierarchical clustering
- ③ Silhouette score
- ④ DBSCAN clustering

{similar kind
of data}



Cluster Ensemble Technique

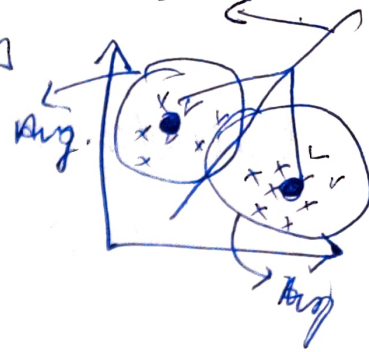
2-3 grp



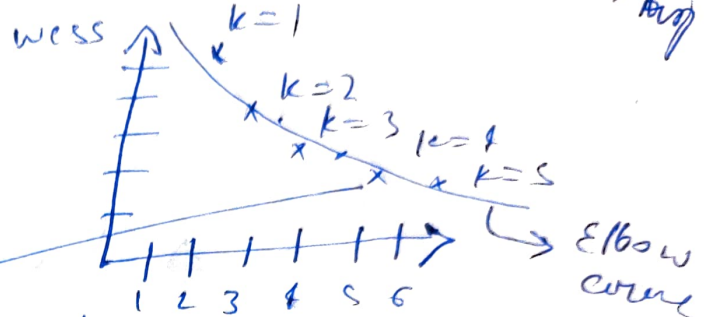
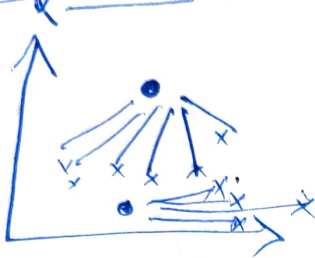
① we try k values \Rightarrow Suitable $k=2$ Euclidean distance
 \hookrightarrow Centroids.

② Initialize k no. of centroids

③ find avg of those



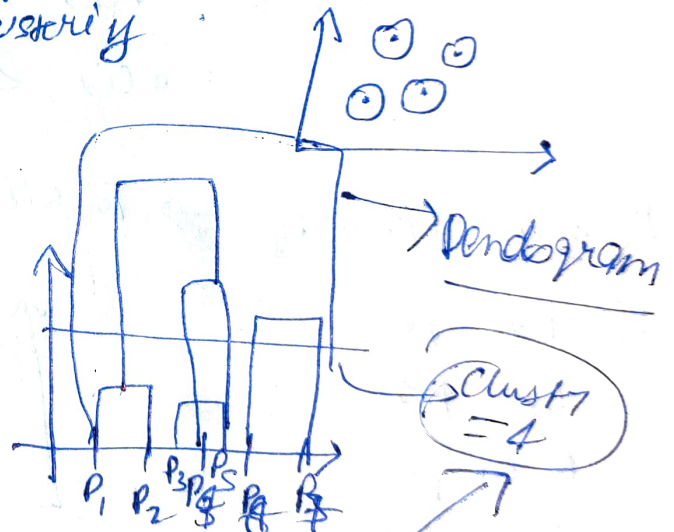
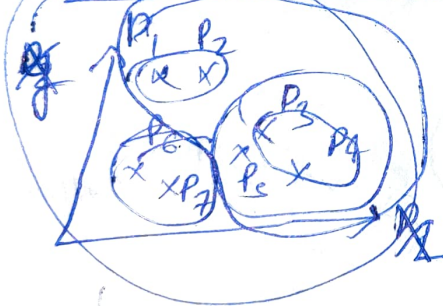
Elbow Method (k value)



Here we find k value where there is abrupt change.

All are the kinds of k -mean clustering

② Hierarchical clustering



You need to find out longest vertical line that has no horizontal line passing through it.

Max time is taken by

k -Mean or Hierarchical
 \hookrightarrow bcz it will keep on creating lots of dendrogram

Dataset is small \rightarrow Hierarchical

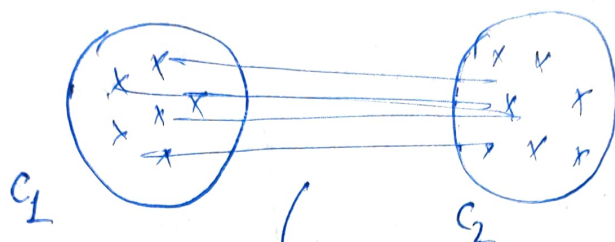
\gg

\gg large \rightarrow k-Mean.

more

Validate clustering models (- to +) \rightarrow good cluster model

\rightarrow Silhouette (clustering).



\Rightarrow Average will be calculated.

$$s(c_i) = \frac{b(i) - a(c_i)}{\max(a(c_i), b(i))}, \text{ if } |c_i| > 1$$

Good cluster $\rightarrow a(c_i) \gg b(c_i)$

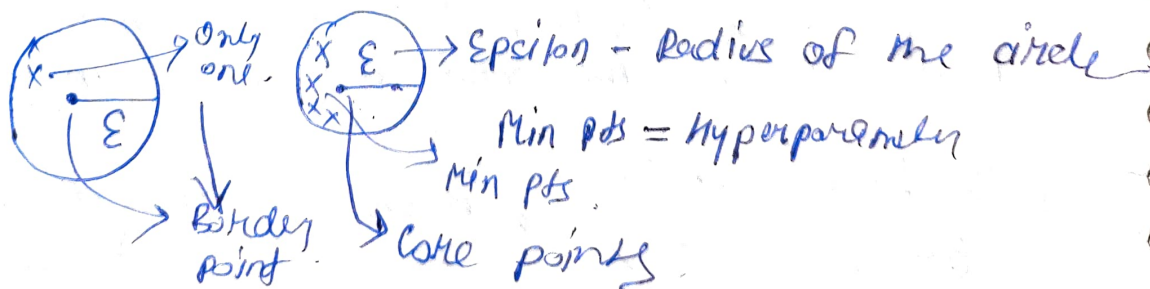
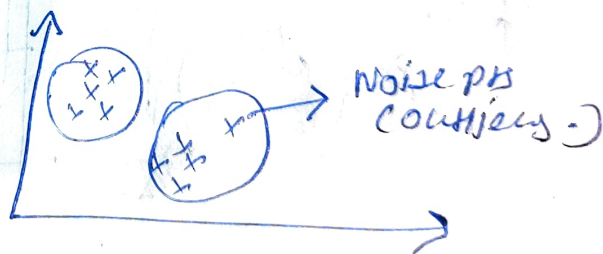
or

$a(c_i) \ll b(c_i) \ll$

DB Scan Clustering

\rightarrow To skip outliers.

- ① Min pts
- ② Core pts
- ③ Border pts
- ④ Noise pts.



\rightarrow If no core or border pts then noise point will be neglected.

K Means Vs

Hierarchical

4280
3/56



Bias & Variance

It is a phenomenon that shows the result of algorithm in favour or against the idea.

Train = 90%
Test = 10% } → Overfitting
 ↓
 Low Bias
 High Variance

Train data → Good - High Bias
 ↓
 Bad - Low Bias

Variance

Refers to the change in model when using different portions of the training or test data.

Test data → Good - Low Var.
 ↓
 Bad - High Var.

SVM, SVR & Xgboost

① Xgboost Classifier

Extreme Gradient Boosting

Boost Model → Pr = 0.5

Dataset

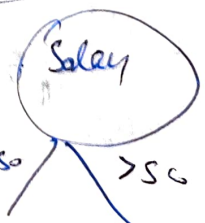
Salary	Credit	Approval	Residual
≤ 50	B	0	-0.5
≤ 50	G	1	0.5
≤ 50	G	1	0.5
> 50	B	0	-0.5
> 50	G	1	0.5
> 50	N	1	0.5
≤ 50	N	0	-0.5

① Calculate a Binary Decision Tree using the feature.

② Calculate Similarity weight

$$= \frac{\sum (\text{Residual})^2}{\sum (P_0(1-P_0) + 1)}$$

$$\text{Sim}_w = \frac{[-0.5, 0.5, 0.5, -0.5] \cdot [-0.5, 0.5, 0.5, -0.5]}{[0.5(1-0.5) + 0.5(1-0.5)] + 0.5(1-0.5) + 0.5(1-0.5)} = 0$$



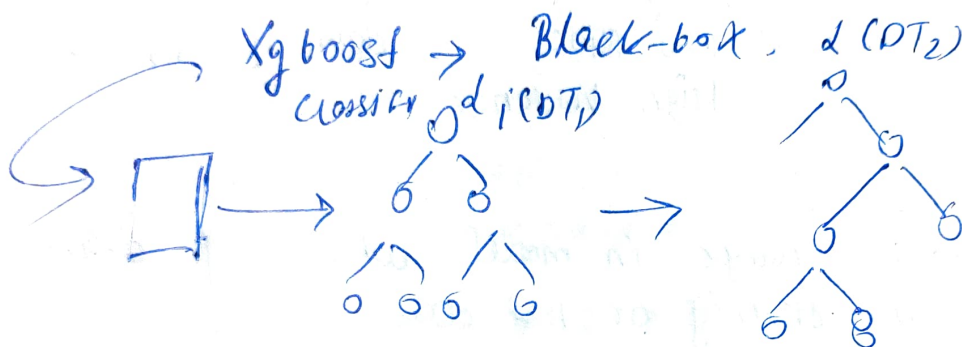
③ Information Gain

$$= 0 + 0.3 - 0.14 = 0.16$$

$$Sim_w = \frac{0.28}{1.75} = \frac{1}{7} = 0.14$$

We basically compare on basis of information gain.

$$O/p \rightarrow \left[C \left[0 + d_1(DT_1) + d_2(DT_2) + d_3(DT_3) + \dots + d_q(DT_q) + \dots + d_n(DT_n) \right] \right]$$

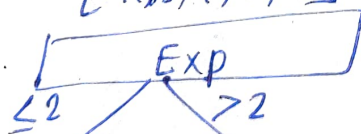


Xg-boost Regressor

Exp	crap	Salary \rightarrow O/P	P_i
2	Y	40K	-11K
2.5	Y	42K	-9
3	N	52K	1
4	N	60K	9
4.5	Y	62K	11

\rightarrow [Base model] \rightarrow SLK

$$[-11, -9, 1, 9, 11] \rightarrow S_m = 1/6$$



GSS

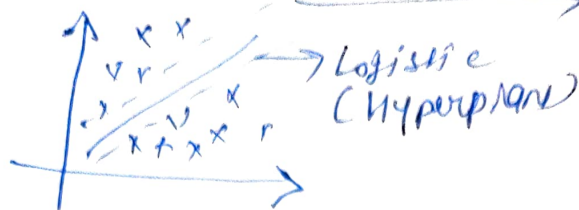
$$[-11] = \frac{12.5 - 65.5}{1+0}$$

$$[-9, 1, 9, 11] = \frac{(-9+1+9+11)^2}{4+1}$$

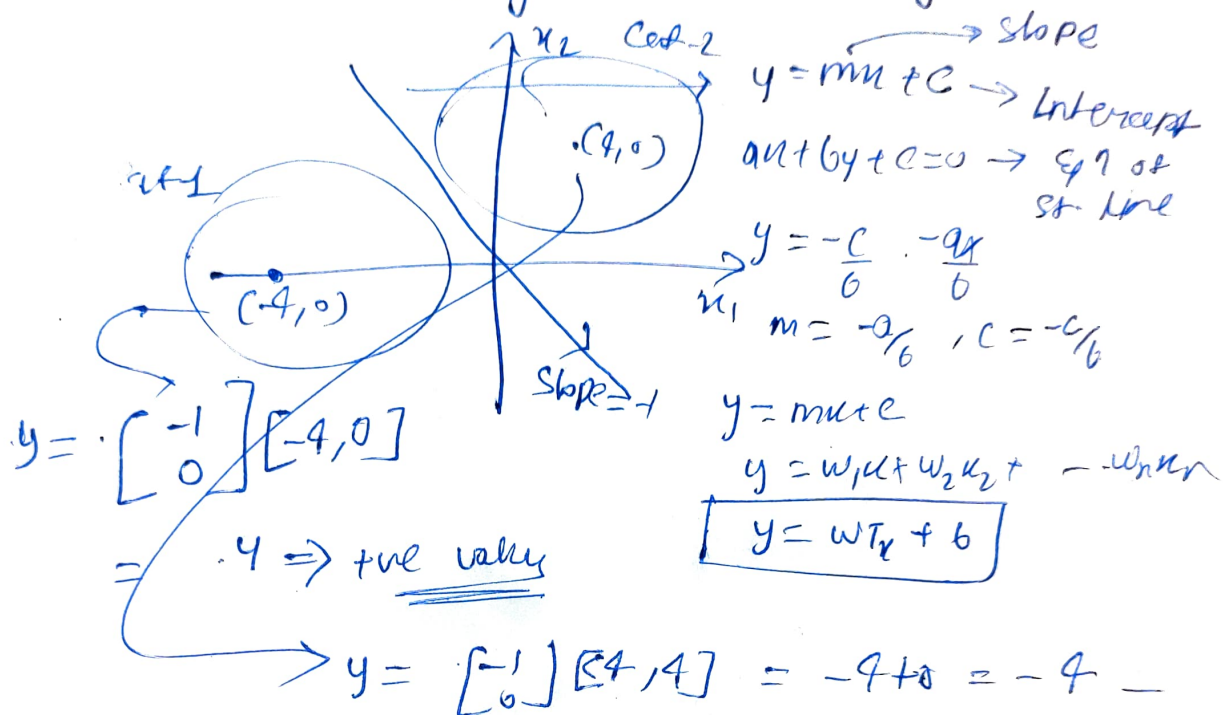
$$Sim\ weight = \frac{\sum (Residual)^2}{No. of residuals + 1} = \frac{144}{5} = 28.8$$

$$Inform\ gain = 65.5 + 28.8 - \frac{1}{6} = 94.34$$

SVM → Similar like logistic regression → Marginal plane



→ In SVM we are going to increase marginal plane.



Now marginal plane is created :-

