

→ Lm. score(X-test, Y-test) - r² score 1

→ mse(Y-test, y-pred) Date.....

[LR]

↳ SKLearn.metrics import mean_squared_error as mse

$$y = \beta_0 + \beta_1 x + \epsilon$$

error term

Intercept

$$\text{residue} = (y_i - \hat{y}_i)$$

sklearn.model_selection
import train_test_split

sklearn.linear_model

import LinearRegression

method of least squares → min of error ↓
 $\sum (y_i - \hat{y}_i)^2$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

sample corr coeff
 $r_{xy} = [\text{sign } g_{b_1}] \sqrt{r^2}$.

$$b_0 = \bar{y} - b_1 \bar{x}$$

* SST → Total sum of squares

SSR → Sum of squares due to regression

SSE → Sum of squares due to error.

$$\hookrightarrow SST = SSR + SSE$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

r^2 → Coeff of determination | $r^2 = \frac{SSR}{SST}$

$$1 = \frac{SSR}{SST} + \frac{SSE}{SST} = r^2 + \epsilon = 1$$

$r^2 = 1 - \epsilon$

Customized
ASIC +
Op-amp

$$r^2 = 1 - \frac{SSE}{SST} \rightarrow 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

Spectrogram

Signal

Specrogram Plot Python

1 - $\frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$
Spiral

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

Date.....

Adjusted R² → used of MLR

why? → As we keep on adding new features, the R² value keeps on ↑ even if the added feature are not correlated to the Target Varible.

Adj R² → It penalizes attributes that are not contributing towards the Target Var.

$$1 - \left[\frac{(1 - R^2)(N - p - 1)}{N - p - 1} \right]$$

↓
No. of predictors
sample size. [No. of samples]

→ Only ↑ when the added Ind var is correlated/ significant and affects the target var.

$$\boxed{\text{Adj } R^2 \leq R^2}$$

— x — x — x —

Four Assumptions of LR.

① Linear Relationship

→ There should be a linear reln b/w x and y.

→ ♦ How to determine if the Assumption is met?

→ Create a scatter plot (x vs y) → points in plot could fall along a straight line.

→ ♦ If Assumption is violated?

→ Apply a non-linear transformation to the variables [log, square root, reciprocal]

Spiral

Date.....

→ Add Another Independent Var to the model.

[if parabolic shape \rightarrow might add x^2 as additional independent variable]

2] Independence

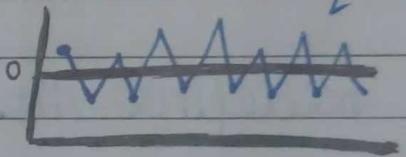
→ Residuals are Independent \rightarrow No correlation/pattern among consecutive residuals.

→ ♦ How to determine if the assumption is met?

→ look at residual time series plot [residue vs time]

→ most of residual autocorrelations should fall within the 95.1% confidence Bands around zero, which are located at about $\pm 2/\sqrt{n}$ (n = sample size)

± 2
over
 \sqrt{n}



→ ♦ If assumption is violated?

→ 1) For +ve serial correlation \rightarrow Consider adding lags of the dependent/independent variable to the model

2) for -ve serial correlation \rightarrow check if none of your variables are overdifferenced.

3) for seasonal correlation \rightarrow adding seasonal dummy variables to the model.

3] Homo scedasticity

→ Residuals have constant variance at every level of x.

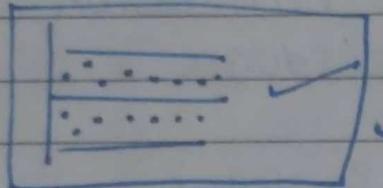
↳ If violated \rightarrow Heteroscedasticity.

↳ ↑ Variance of the

Coeff \rightarrow declares the term significant

Spiral when it is not.

- How to check if assumption is met.
- fitted value VS Residual plot [scatter plot]



* every level of x the residual gets more apart (\uparrow variance)

- How to fix, if violated.

① Transform the dependent variable. → log of dependent var.

② Redefine the dependent variable → uses
→ use rate, rather than raw value.

③ use weighted regression → assigns a weight to each data point based on the variance of its fitted value.

→ gives small weight to p if that have higher variance.

4] NORMALITY

→ residuals are normally distributed.

① How to check.



① Using Q-Q plots

Date.....

if the points on the plot follow/ form a straight diagonal line, then assumption is met.

if violated

- 1) Verify that outliers aren't having a huge impact on the distribution.
- 2) non-linear transformation on variables.

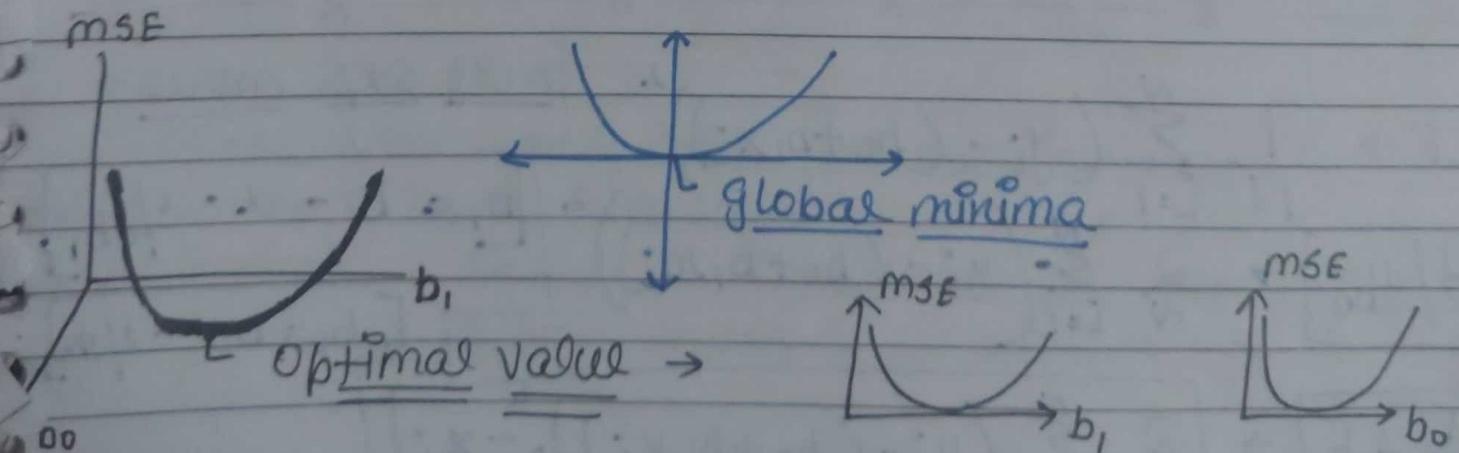
x x x
Cost funcⁿ and gradient descent.

↳ find optimal values of β_0 and β_1 .

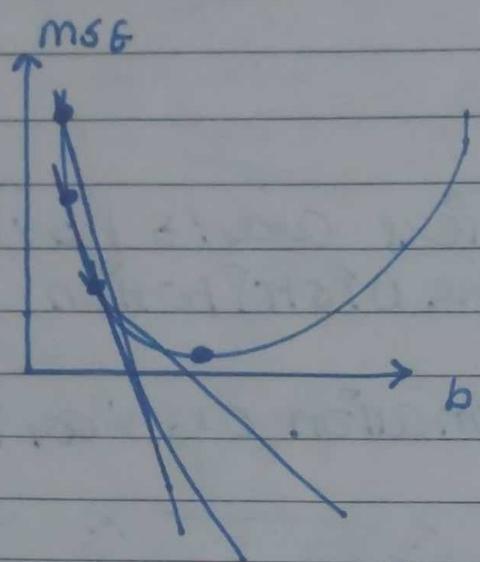
minimum error \hookrightarrow global minima [values of β_0 and β_1 , where error is minimum \rightarrow global minima].

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

$$\text{cost func}^n \rightarrow \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \rightarrow \frac{1}{N} \sum_{i=1}^N (y_i - (b_0 + b_1 x_i))^2$$



• Gradient Descent → Optimization Algo that helps our cost function to reach optimal point (values with min error).



① -

Slope at every point → To get the direction where we want to move.

→ changing value of b_0 and b_1 , so that we can get values with min error → optimal value.

* No. of steps / the method → Learning Rate.

→ High LR → overshoot → miss the optimal value.

→ small LR → more steps But will reach the optimal value.

② To take the next step → Learning Rate

 (α)

$$CF \rightarrow \frac{1}{N} \sum_{i=1}^N (y_i - (b_0 + b_1 x_i))^2$$

Next step move →

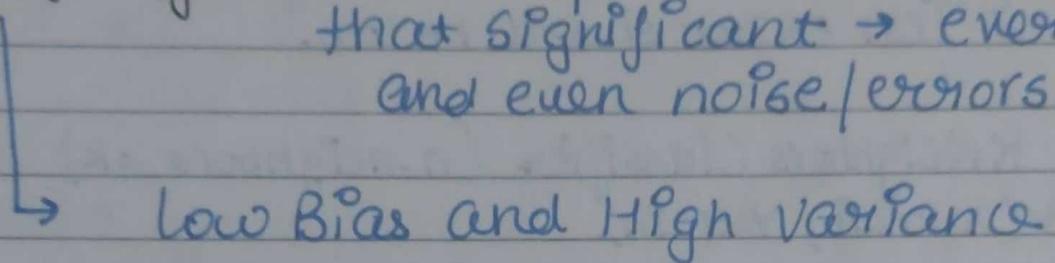
$$\boxed{\frac{d}{db_0}} \rightarrow \frac{2}{N} \sum_{i=1}^N (y_i - (b_0 + b_1 x_i))^2 \quad \boxed{b_1 = b_1 - \alpha \cdot \frac{d}{db_1}}$$

$$\boxed{\frac{d}{db_1}} \rightarrow \boxed{\frac{2}{N} \sum_{i=1}^N (y_i - (b_0 + b_1 x_i))^2 \cdot -x_i}$$

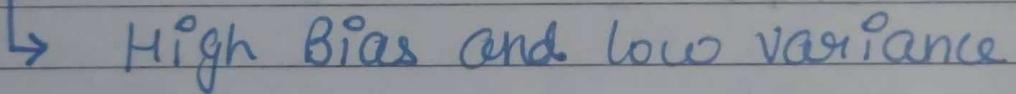
$$\boxed{b_0 = b_0 - \alpha \cdot \frac{d}{db_0}}$$

- High Training accuracy \rightarrow low Bias
(Train data error)
- low Test accuracy \rightarrow High Variance
(Test data error)

Overfitting \rightarrow model captures noise too - makes that significant \rightarrow every pattern and even noise/errors.



Underfitting \rightarrow model fails to learn from Training data set \rightarrow fails to capture every pattern or significant patterns.



Best fit \rightarrow low Bias and low variance.

\rightarrow X — X — X — X —

KNN

- measure distance of points from the test data
- sorts them in ↑ measure
- depending on no. of K, considers K data points
- majority class from those AP is assigned to Test Data.

→ Distance Based learning Algo → Feature Scaling Date.....

① Choose value of K.

① Should be odd → Not multiple of no. of classes.

② Error Curves → (curves of error of Training and Testing Data)
Test error should be low.

Knn = K Neighbors Classifier (n-neighbors = K)

Knn.fit(X-train, Y-train).

y-pred1 = Knn.predict(X-train)

error1.append(np.mean(y-train != y-pred1))

y-pred2 = Knn.predict(X-test)

error2.append(np.mean(y-test != y-pred2))

plt.plot(error1). [Distance → $\sqrt{(x_2-x_1)^2 + \dots}$].
plt.plot(error2).

• Performance Metrics

① FPR ② FNR ③ Accuracy ④ Recall

⑤ Precision ⑥ F-Beta

⑦ Cohen's Kappa Stats

• Confusion Matrix :

TP Pred 1
FP Pred 0
Actual Result
of Both Actual
and Predict

		1	0	→ Actual val
Pred values	1	TP	FP	→ Type 1 error (FPR)
	0	FN	TN	→ Type 2 error (FNR)

Date.....

$$FPR = \frac{FP}{TN+FP} \text{ (negative)}$$

$$FNR \rightarrow \frac{FN}{TP+FN} \text{ (positive)}$$

$$TPR \rightarrow \frac{TP}{TP+FN}$$

$$TNR \rightarrow \frac{TN}{TN+FP}$$

★ → TPR and TNR Should be high.
FPR and FNR Should be Low.

• Accuracy = $\frac{TP+TN}{N}$

• Precision → $\frac{\text{True +ve}}{\text{pred +ve}} \rightarrow \frac{TP}{TP+FP}$

→ Out of Total predicted +ve results
how many results were actual +ves.

• Recall → $\frac{\text{True +ve}}{\text{Actual +ve}} \rightarrow \frac{TP}{TP+FN}$

→ Out of Total actual +ve, how many
+ve did we predict correctly.

* when Prec and Recall Both are Imp, For that
we use F-Beta

$$\hookrightarrow \frac{(1+\beta^2) \left[P \times R \right]}{\beta^2 (P+R)} \left[\frac{(1+\beta^2) \left[P \times R \right]}{\beta^2 (P+R)} \right]$$

Date.....

when $\beta = 1 \rightarrow \underline{\underline{f_1 \text{ score}}}$

$$= \frac{2(P \times R)}{P+R} + \boxed{\frac{2PR}{P+R}}$$

* Cohen's Kappa stats \Rightarrow can handle Both multi-class and imbalanced class problems.

$$\boxed{K} \rightarrow \frac{(p_o - p_e)}{(1 - p_e)}$$

$p_o \rightarrow$ Observed agreement

$p_e \rightarrow$ Expected agreement.

- $K < 0 \rightarrow$ No agg
- $0 \leq K < .2 \rightarrow$ slight
- $.21 < K < .40 \rightarrow$ fair
- $.41 < K < .60 \rightarrow$ moderate
- $.61 < K < .80 \rightarrow$ substantial
- $.81 < K < .99 \rightarrow$ near-perfect

$K = 1 \rightarrow$ perfect

$$\boxed{p_o} \rightarrow \frac{[(\text{Both said Yes}) + (\text{Both said No})]}{\text{Total Ratings}}$$

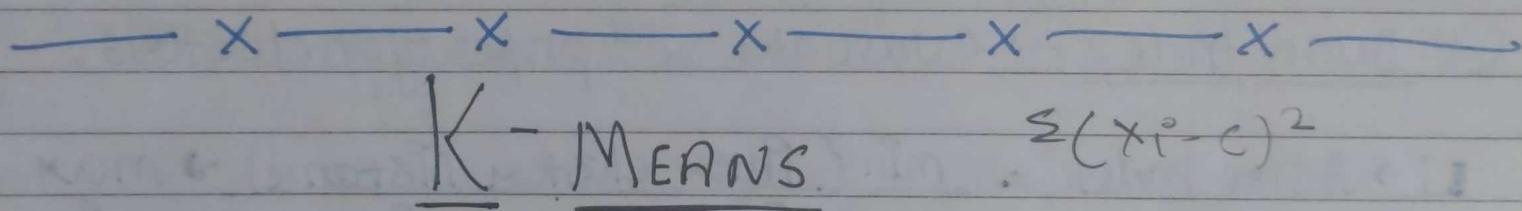
$$p_e \rightarrow P(\text{Yes}) \rightarrow \left[\frac{\text{Rater 1(Yes)}}{\text{Total}} \times \frac{\text{Rater 2(Yes)}}{\text{Total}} \right]$$

$$P(\text{No}) \rightarrow \left[\frac{\text{Rater 1(No)}}{\text{Total}} \times \frac{\text{Rater 2(No)}}{\text{Total}} \right]$$

$$\boxed{p_e \rightarrow P(\text{Yes}) + P(\text{No})}$$

Date.....

- Eager Learners → Given set of Train Data → Create a model Before receiving test/new Data.
- Lazy learners → Simply stores the train data → waits until it is given a Test Data. → No model creation.
 - Less time in Training
But more time in predicting / Testing
 - No training Stage. → all the work is done during test stage.

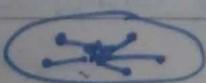


- Clustering → Dividing dataset into groups
 - Here we do not have a Target value - we look at Data and try to group similar observation and form diff grp/ clusters.

Property 1 → All the DP in a cluster should be similar to each other.

Property 2 → DP from different clusters should be as different as possible.

Intra-cluster Dist → Distance of a point from centroid in a same cluster



Inter-cluster Dist → Distance b/w centroids of 2 diff clusters.

Evaluation Metrics

① Inertia → evaluates 1st property of clusters.

① → sum of ~~all~~ intra-cluster distance

② we want less value of Inertia → so that DP in a cluster are close / similar to each other.

→ minimize the Intra-cluster distance ⇒ Inertia should always be low.

② Dunn Index → evaluates 2nd property of clusters.

$$\square \rightarrow \text{Dunn Index} = \frac{\min(\text{Inter cluster distance})}{\max(\text{Intra cluster distance})} \rightarrow \max \min(\text{Intra cluster distance}) \rightarrow \min$$

→ helps in separating diff clusters so that DP from diff clusters are very far from each other.

→ we want max value of Dunn Index (maximize Dunn Index)

- Num should be max → (min of inter-cluster dist)
 - ↳ so that dist b/w even the closest clusters are more (min dist b/w clusters should be max).

- Denom should be min → max of intra-cluster dist
 - ↳ so that max dist b/w points and centroid should be min.

Date.....

* $\rightarrow R_L \rightarrow \infty$

Working \rightarrow min the Distance of the points within a cluster with their centroid.

\rightarrow Calc Distance \rightarrow Assign a point to a cluster.

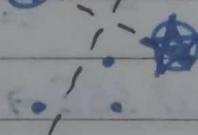
Eg \rightarrow

..
..
.. ..

Step 1 \rightarrow Choose the no. of clusters K

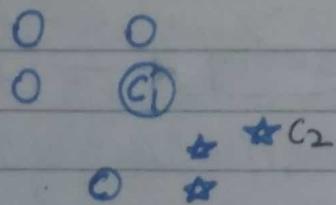
Step 2 \rightarrow Select 'K' random points from the data as centroids.

\rightarrow .. . , — way of assigning points graphically



Step 3 \rightarrow Assign all the points to the closest cluster centroid

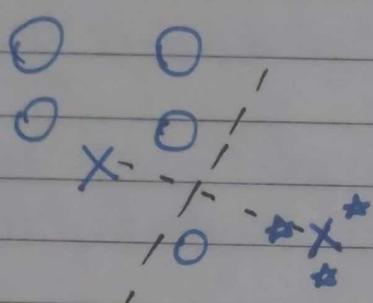
$\rightarrow \rightarrow \rightarrow$ calc Distance of O_P from C_1 and C_2 and assign O_P to that centroid where the Distance is minimum.



→ Calculate mean of each cluster = new centroid
 $\hookrightarrow (x_1 + x_2 + \dots + x_n)/n, (y_1 + y_2 + \dots + y_n)/n \dots$ Date.....

Step 4 → Recompute the centroids of a newly formed clusters.

↳ Adjust centroids of the new clusters so that they become center of gravity of given cluster.



[$x \rightarrow$ new centroid]

Step 5 → Repeat Step 3 and Step 4

Select model (k) with least SSE

$$SSE_1 \rightarrow \sum_{i=0}^N (x_i - c_1)^2$$

$$SSE_2 = \sum_{i=0}^N (x_i - c_2)^2$$

Total SSE $\rightarrow SSE_1 + SSE_2 \dots + SSE_k$

* When to Stop

① Centroids of newly formed cluster do not change

② Points remain in the same clusters
 \rightarrow No new points are being introduced in newly formed clusters.

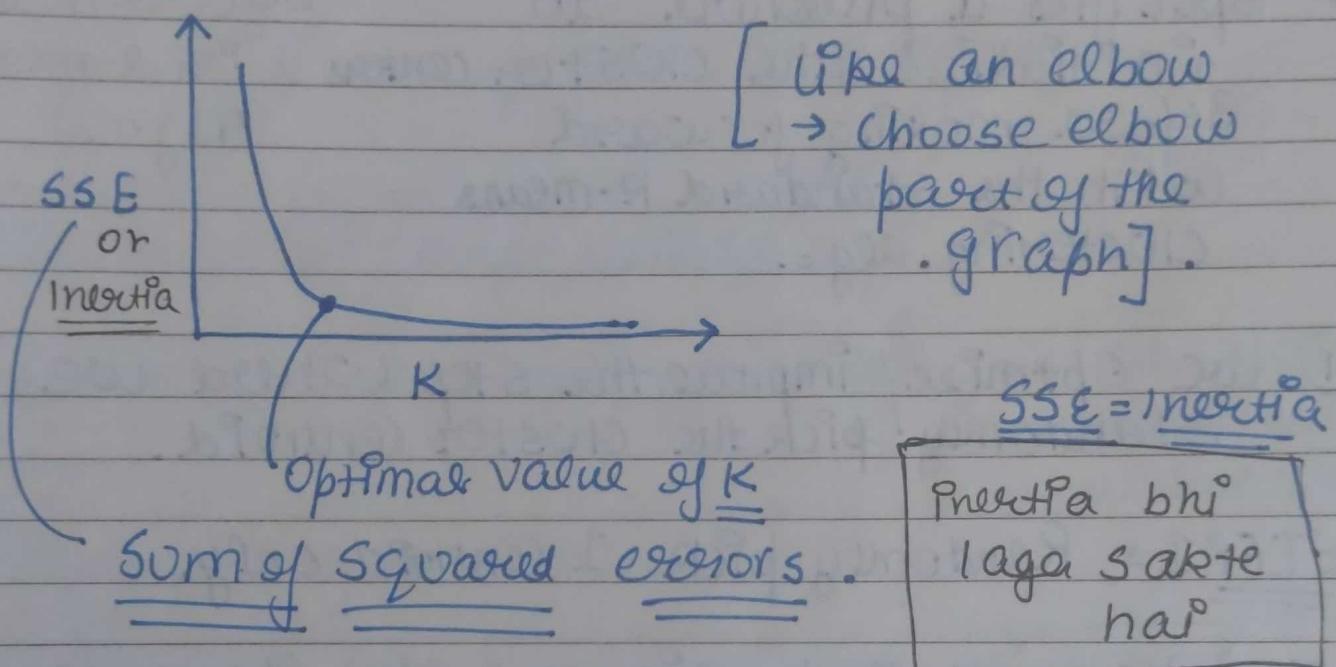
IMP. → Classification via clustering ← IMP.

• Test data placed in cluster. Its class = majority of that cluster
 ↳ can also give probab ($80\% \rightarrow A | 20\% \rightarrow B$), Spiral Class

Date.....

How to find Best Value of K?

- ① Draw a elbow curve
- ② Choose that 'K' where the graph starts to converge into constant or reach to 0.



from `sklearn.cluster import KMeans`

`Rmeans_obj = KMeans(n_clusters=2, init='K-means++')`

```
task [kmeans_obj.fit(data_scaled)  
pred = Rmeans_.predict(data_scaled)]  
Rmeans_obj.inertia_
```

elbow - curve ($SSE = []$)

for cluster in range(1, 20):

`K = KMeans(n_clusters=cluster, init='K-means++')`

`K_.fit(data_scaled)`

`SS_E.append(K_.inertia_)`

Spiral

Date.....

Df = pd.DataFrame($\{ 'K': \text{range}(1, 20), 'SSE': \text{SSE} \}$).
plt.plot(Df['K'], Df['SSE'], marker='o').

① K-means++ / k-means++ \rightarrow Initializer of centroids
 \rightarrow Specifies a procedure to Before Proceeding
initialize the cluster centers with K-means
Before moving forward Algo.
with the standard K-means
clustering algo.

① we Optimize/ Improve the step where we randomly pick the cluster centroid.

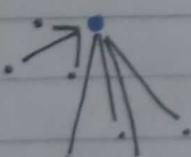
STEP 1 \Rightarrow Randomly pick 1 centroid only.

STEP 2 \Rightarrow Compute the Distance $D(x)$ of each $DP(x)$ from the cluster center that was chosen.

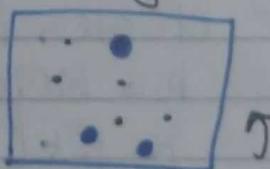
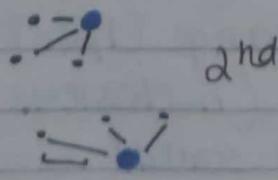
STEP -3 \Rightarrow Choose new cluster center from the DP with the probab of $P_C \propto D(x)^2$.

STEP -4 \Rightarrow Repeat Step 2 and 3 until K-cluster centers are chosen.

Then 2, 5, for
calc DP of st
to Both K
choose the
min one and
then select it



\rightarrow new centroid $\rightarrow (D(x_1))^2$ is
farthest



3rd \rightarrow Distance of each point from closest centroid and the point with largest $(D(x))^2 \rightarrow$ centroid spiral

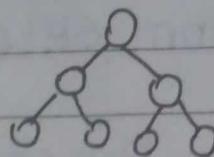
overfitting \leftrightarrow low Bias and High Variance

DT

Date.....

Binary split

f_1, f_2, f_3, \dots output



Binary classification

- 1) Entropy \rightarrow Measure of the purity of split
 \rightarrow Helps in deciding which feature should be selected for splitting.

Value ranges from $0 \rightarrow 1$ | Low value - Better

- Pure split \rightarrow $\{1\text{ Yes} / 0\text{ No}\}$ OR $\{0\text{ Yes} / 1\text{ No}\}$
 \downarrow
 $\text{Entropy} = 0$ (usually leaf node).

- Impure split \rightarrow $\{1\text{ Yes} / 1\text{ No}\} \rightarrow \text{Entropy} = 1$

Imp

- It only counts for single node But in splitting we need to calculate for other nodes also. As they are part of tree so we use Information Gain which includes entropy.

$$E \rightarrow H(s) = -P_{(+)} \log_2(P_+) - P_{(-)} \log_2(P_{-})$$

$\rightarrow P_{(+)} = \text{\% of +ve class}$
(Probab)

$\rightarrow P_{(-)} = \text{Probab of -ve class}$.
[3-Yes | 2-No]

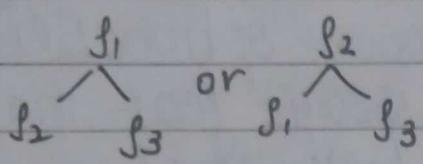
$$= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$

$$= \underline{\underline{0.79 \text{ bits}}} \rightarrow \text{Not good}$$

Spiral

Date.....

- Information Gain \rightarrow avg of all Entropies of a tree to see which split is better.

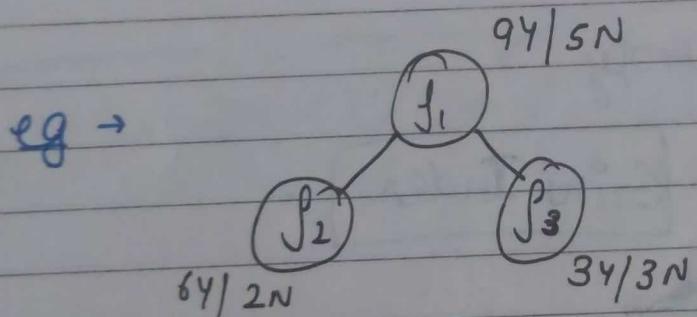


on which feature? using purity of split of each node and using their average
↑ Imp

$$\text{Grain}(S, A) = H(S) - \sum_{\substack{\text{VAL} \\ \text{Root}}} \frac{|S_v|}{|S|} H(S_v).$$

$$\rightarrow E_S = \sum_{\substack{\text{VAL} \\ \text{Root}}} \frac{|S_v|}{|S|} E_{Sv}$$

subset after splitting.



$$\begin{aligned} S &\rightarrow 14 \\ S_{f_2} &\rightarrow S_{v_1} \rightarrow 8 \\ S_{f_3} &\rightarrow S_{v_2} \rightarrow 6 \end{aligned}$$

$$\rightarrow E_S = \sum_{\substack{\text{VAL} \\ \text{Root}}} \frac{|S_v|}{|S|} E_{Sv} \quad \left| \begin{array}{l} E_{f_1} = 0.94 \\ E_{f_2} = 0.81 \\ E_{f_3} = 1 \end{array} \right.$$

$$\rightarrow 0.94 - \frac{8}{14} \cdot 0.81 - \frac{6}{14} \cdot 1 = \underline{\underline{0.049}}$$

* Higher the value - Better it is. - That way it's better for splitting.

Date.....

* Gini Impurity [Similar to Entropy But Better].

$$\rightarrow 1 - \sum_{i=1}^n (P_i)^2$$

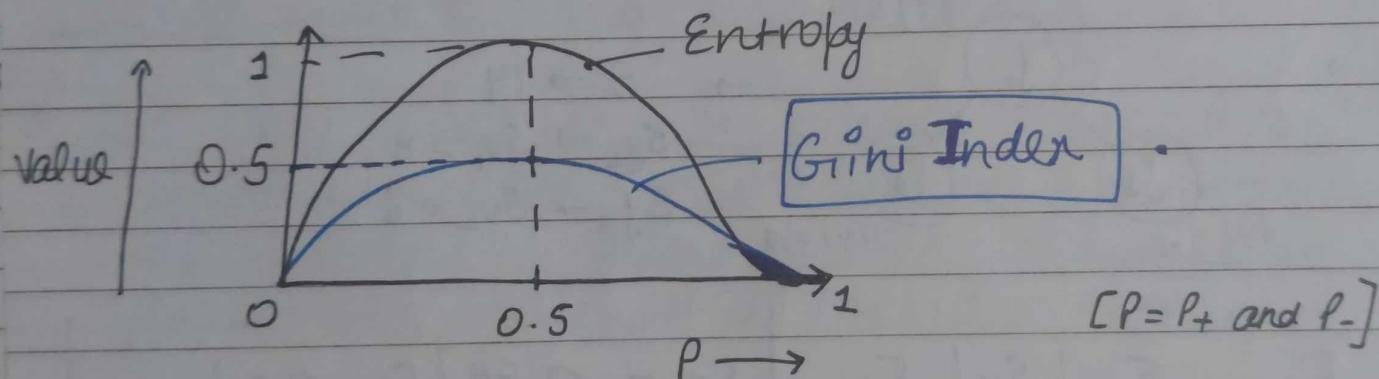
$$= 1 - [P(+)^2 + P(-)^2]$$

Amount of Probability of a specific feature that was classified incorrectly when selected at Random

\rightarrow Computationally Efficient

If I have 34 and 3N $\rightarrow E = 1$

But Gini Index would be 0.5



\rightarrow Shorter period of Time for Execution.
(log takes time).

Imp

- Then we can use Information gain by using entropy or Gini Impurity.

- Gini - CART Algo
- IG \rightarrow ID₃ algo.

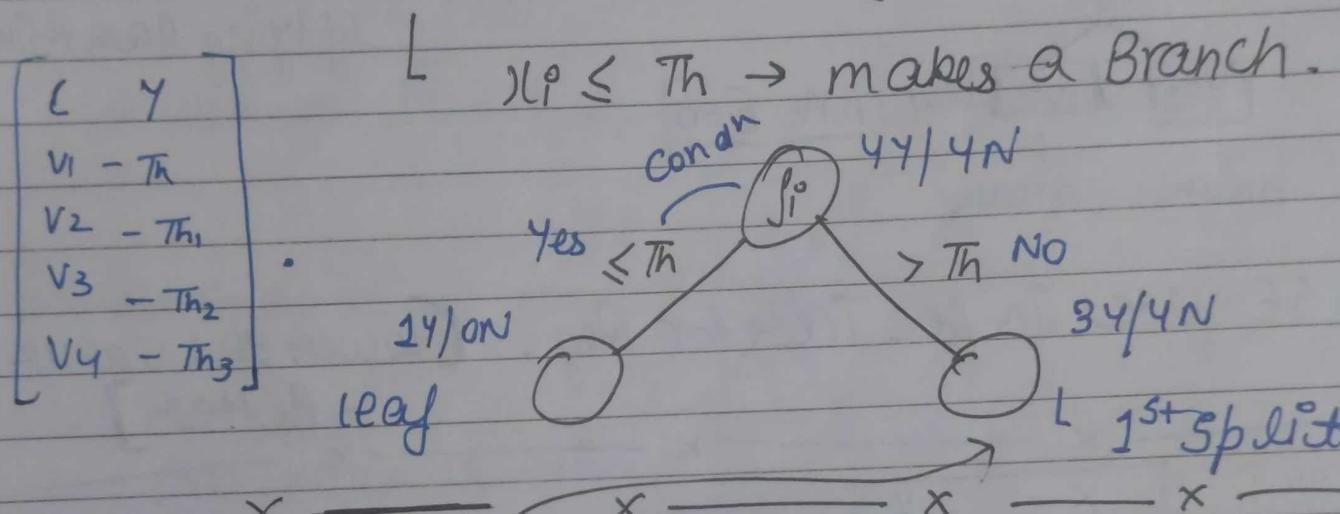
.. DT for Numerical feat - Date.....

→ R → DT will ↗

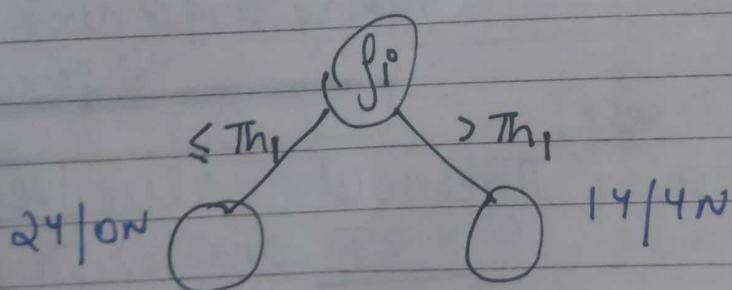
→ ① Sort all the values of that feat

② Take Threshold Value. (Th).

2.1 select a value (every value go - through)



• for 2nd Split (2nd way of Split).



check / split
for every threshold
happens

• which is Better (with Th or with Th_1)?

↳ Entropy and If

↳ Better If will be taken.

* Takes Time

• If of every feature is calculated wrt to the output variable in normal data.

Ensemble learning Date.....

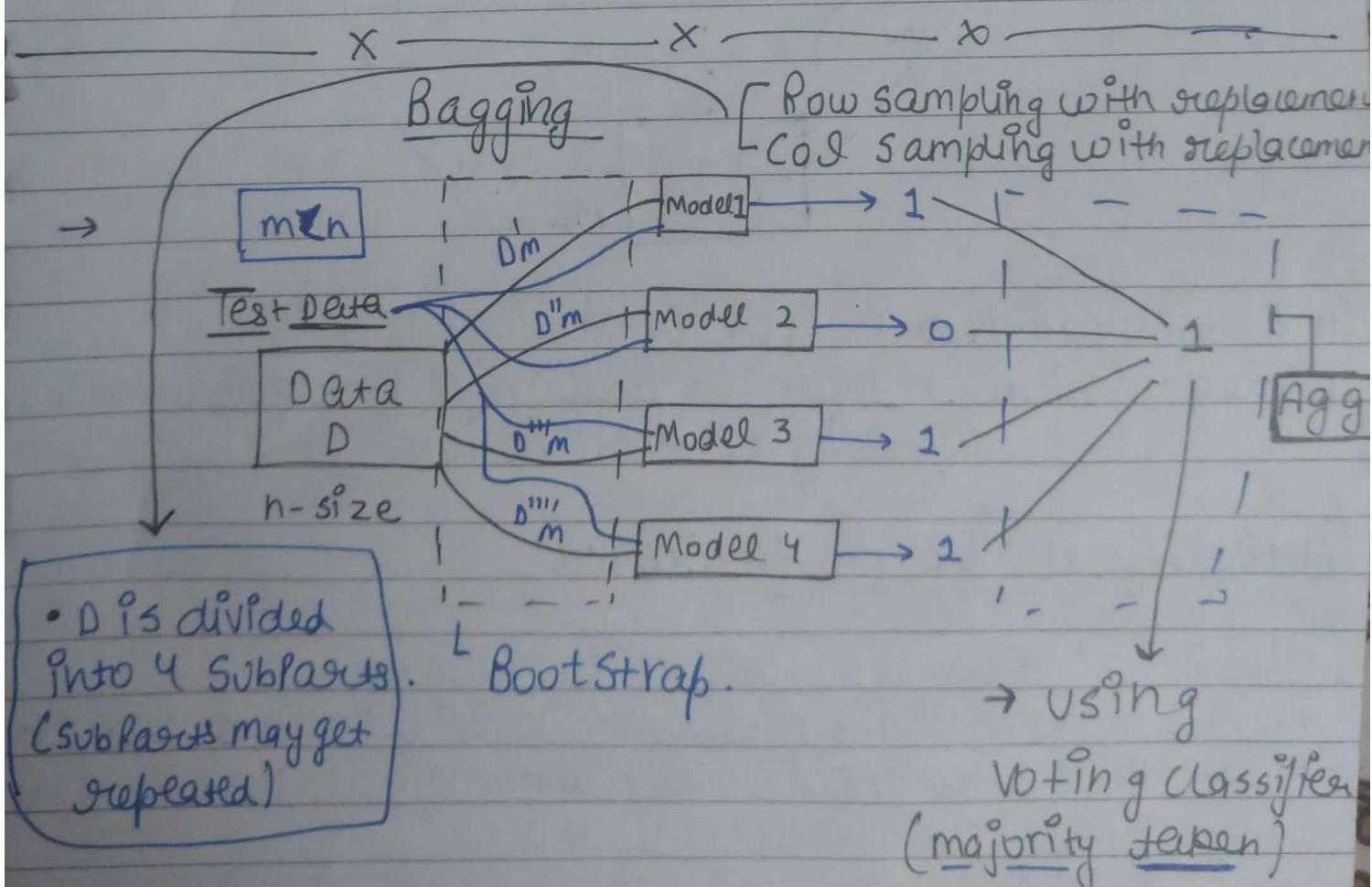
→ Combining Multiple Models.

① Bagging (Bootstrap - Divide)
(Aggregation - Combine)
(Algorithm).

→ Random Forest

② Boosting

- AdaBoost
- Gradient Boosting] Advance.
- XgBoost



Random Forest

Date.....

→ Bagging Tech

- RS → Row sampling with replacement
- CS → Col sampling with replacement

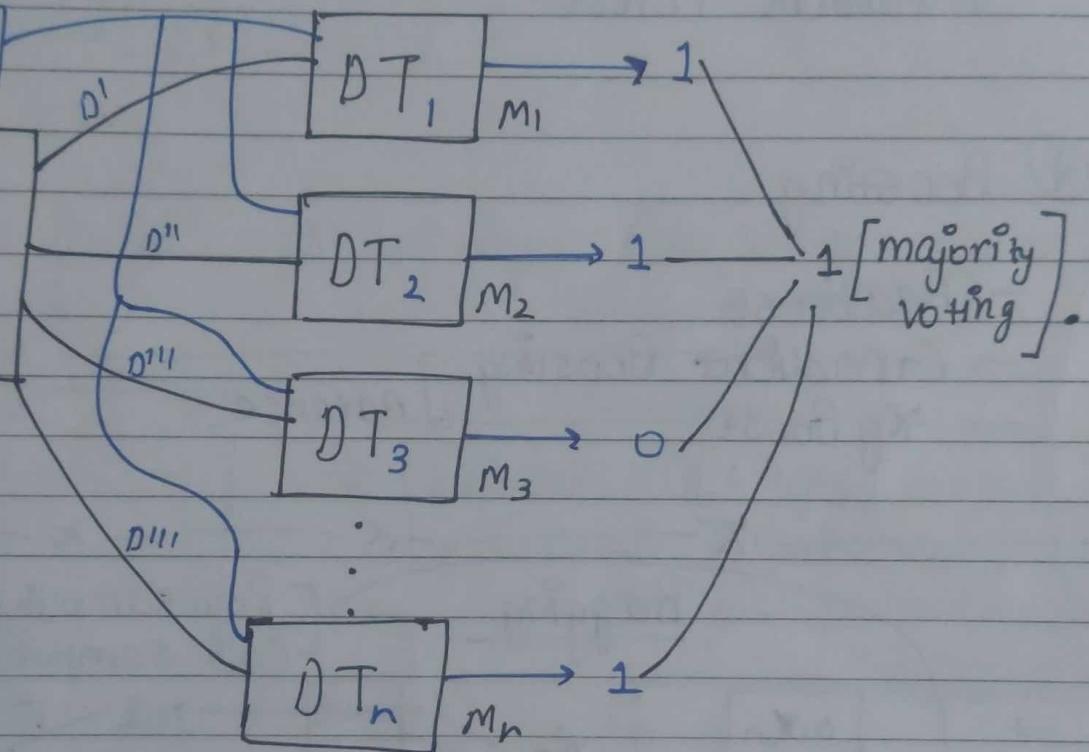
$$RS + CS \rightarrow D'$$

Test Data

Data

D

n



In Classification → Majority Voting.

In Regression → mean or median of 'n' output values depending on distribution.

Advantage ① Low Variance

- In DT → low Bias and High Variance
- + When multiple DTs are used in RF with RS and CS
 - ↳ get TRAINED SPECIFICALLY and majority vote converts High Var to low Var.

② Change of m rows (m < n) will not affect the model that much.