

ML ALGO^M

PAGE NO.:

DATE: 11

DAY-1 INTRODUCTION TO ML ALGO

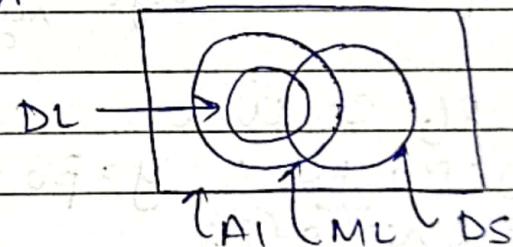
1. AI VS ML VS DS VS DL

- AI appⁿ: able to do its task without any human intervention

ex: Netflix → Recommendation system
Amazon, self driving cars.

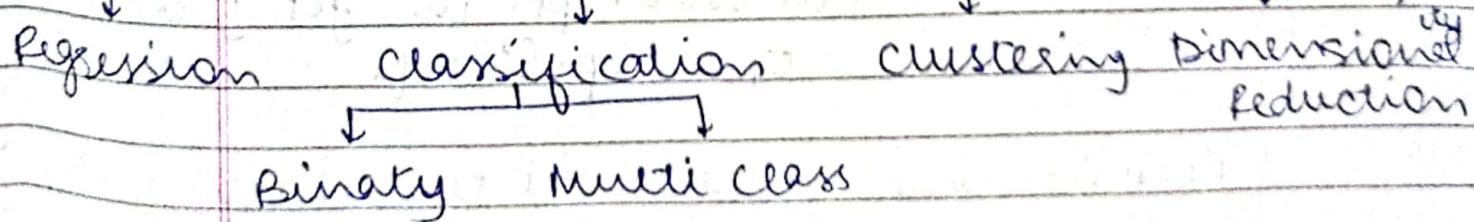
- ML: subset of ai; state tool to analyse; visualise; prediction; forecasting.

- DL: mimic human brain; mimic human brain



Machine Learning

Supervised L. & Unsupervised L.



* Supervised Algoms

linear regression, Ridge & lasso,
logistic regression, decision tree, Adaboost,
Random Forest, Gradient Boosting,

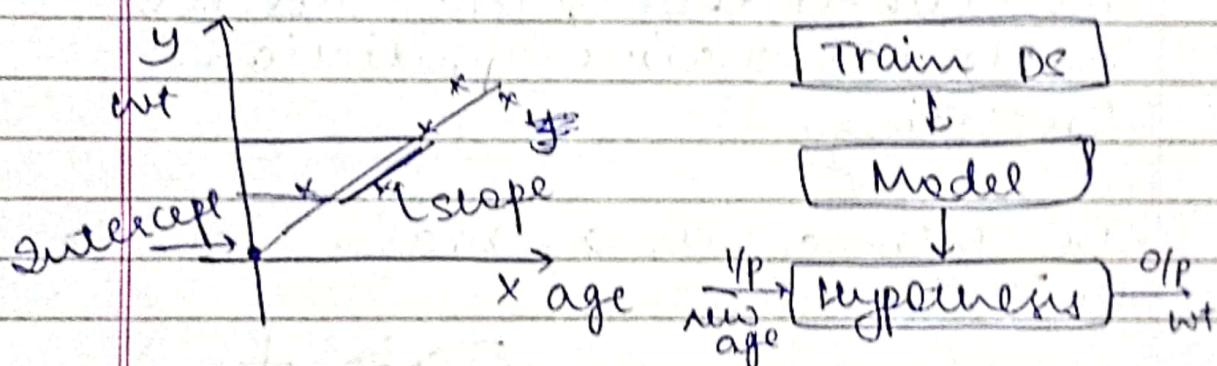
Machine Learning

XGBoost, Naive Bayes, SVM, KNN

x. Unsupervised learning

K-means, DBScan, Hierarchical clustering,
K-nearest neighbour cluster, PCA, LDA

1. Linear Regression



Eqn of st. line

$$y = mx + c ; \quad y = \theta_0 + \theta_1 x ;$$

Hypothesis (\hat{y})

$$\hat{y} = \theta_0 + \theta_1 x$$

θ_0 = Intercept

When $x=0$, $y=\theta_0$

at what pt. we get y axis

θ_1 = slope % coeff.

Best-fit line: When all the data pts
are at min distance from the
prediction line.

Cost F_m:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \rightarrow \text{Squared Error Function}$$

L_m: average deviation per point or similar calc

what we need to solve?

$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$$

$$* h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$

$$\text{If } \theta_0 = 0$$

Intercept is
passing through
origin.

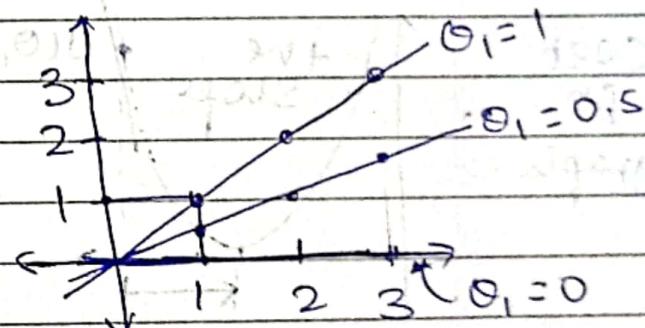
$$h_{\theta}(x) = \theta_1 \cdot x$$

$$h_{\theta}(x) = \theta_1 \cdot x$$

$$(1, 1) \quad \text{plus}$$

$$(2, 2)$$

$$(3, 3)$$



$$\theta_1 = 1$$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$$

$$= \frac{1}{2m} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$$J(\theta_1) = 0$$

$$\theta_1 = 0.5$$

$$J(\theta_1) = \frac{1}{2m} [(.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

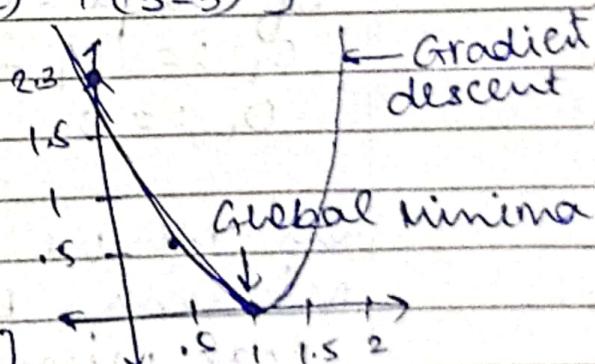
$$= \frac{1}{2m} [.25 + 1 + 2.25]$$

Cost $J(\theta_1)$

$$= 0.58$$

$$\theta_1 = 0$$

$$J(\theta_1) = \frac{1}{2m} [(0-1)^2 + (0-2)^2 + (0-3)^2] \approx 2.3$$



Convergence Algorthm

Repeat until convergence

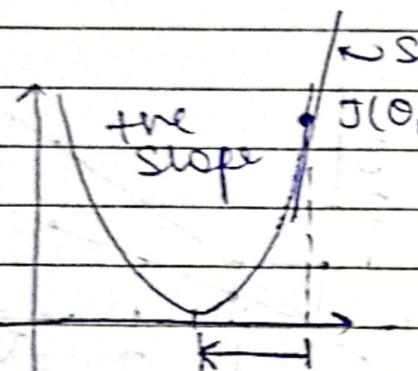
{ learning Rate

$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

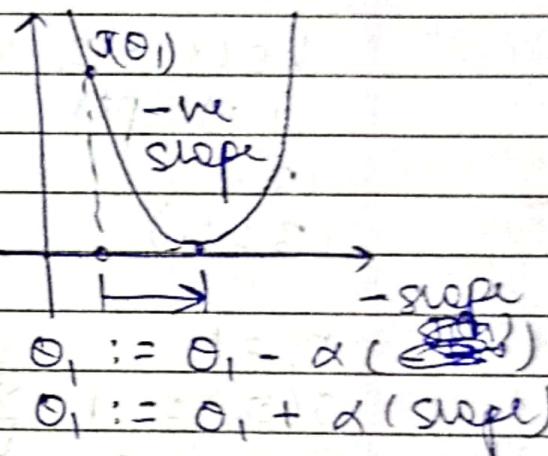
}

derivative of J

slope.

Cost fn
graph.

$$\theta_1 := \theta_1 - \alpha (+ve)$$



$$\theta_1 := \theta_1 - \alpha (-ve)$$

$$\theta_1 := \theta_1 + \alpha (slope)$$

if slope = 0

$$\theta_1 := \theta_1 - \alpha (0)$$

$$\theta_1 := \theta_1$$

In ML we don't have any local minima; But in DL we may have local minima and we are stuck.

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i)^2$$

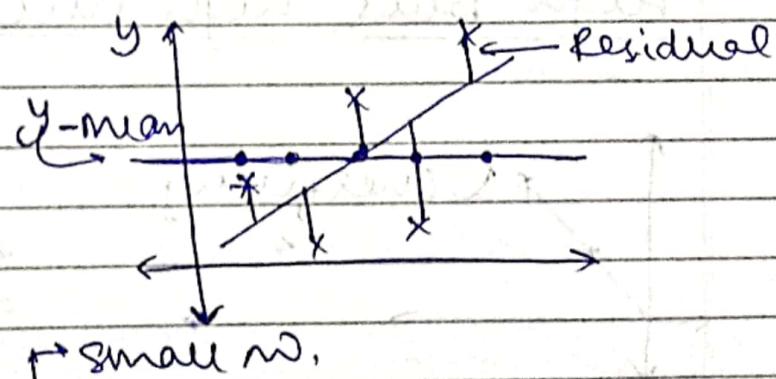
$$j=0 \Rightarrow \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i)$$

$$j=1 \Rightarrow \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) \cdot x^i.$$

Performance Matrix

↳ R^2 and Adjusted R^2 .

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{total}}} = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2}$$



$$= 1 - \frac{(\text{low})}{(\text{high})} = \text{Big NO} = 90\%$$

If we add a feature then our R^2 will inc. But if we add a completely independent feature that is not correlated then also the R^2 inc. This is a problem.

$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(N-1)}{N-P-1}$$

P: features or predictors

N: no. of data points

DAY: 2 Ridge, Lasso logistic regression

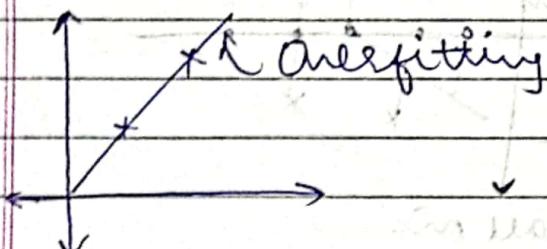
Ridge & Lasso regression

Overfitting: Model performs well with training data but fails to perform well in ~~training~~ with test data.

Low Bias, High Variance
Model

Underfitting: Model accuracy is bad with training data and model accuracy is also bad with testing data.

High Bias, Low Variance



$$h_0(x) = \hat{y}$$

~~$J(\theta_0) = 0$~~

$$= \frac{1}{2m} \sum_{i=1}^m (h_0(x^i) - y^i)^2$$

$$\text{Let, } \lambda = 1 \text{ and } J(\theta_0) = \frac{1}{2m} \sum_{i=1}^m (y - \hat{y})^2$$

In ridge (L2 regularization)

$$J(\theta_0) + \lambda (\text{slope})^2$$

$$J(\theta_0) = \theta_0 + 1(2)^2$$

$$= 4$$

$$= \sum_{i=1}^m (y - \hat{y})^2 + \lambda (\text{slope})^2$$

small value

$$= (\text{small value}) + \lambda (\text{slope})^2$$

$$= \frac{1}{3} + 1(1.5)^2$$

- Ridge helps in creating generalized model. Helps avoid overfitting.

$\lambda \rightarrow$ hyperparameter

↳ steepness of the slope

LASSO (L1 Regularization)

$$= (\hat{y} - y)^2 + \lambda |\text{slope}|$$

↳ feature selection

$$h_{\theta}(x) = \hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$|\text{slope}| = |\theta_0 + \theta_1 + \dots + \theta_n|$$

(i) Prevents Overfitting

(ii) Feature Selection

Ridge Regression (L2 Reg.)

$$\text{C.F.} = (h_{\theta}(x))^2 + \lambda (\text{slope})^2$$

• Prevents Overfitting

Lasso Regression (L1 Reg.)

$$\text{C.F.} = (h_{\theta}(x))^2 + \lambda |\text{slope}|$$

• Prevents Overfitting

• Feature Selection

#. Assumptions of linear regression

(i) If all features are in normal or Gaussian distribution, our model will be trained well.

(ii) Standardization (Scaling the data by z-score) (wherever there is gradient descent it's good to do standardization.)

(iii) It works w.r.t linearity: If our data is linearly separable then it will be able to give good ans.

(iv) Multicollinearity

Logistic Regression (Classification)

- Works well with binary classification
- Can also use for multi class.

Ex: Predict P/F based on no. of hrs study.

Q. Why can't we use linear regression for classification?

- Because of an outlier our line is completely shifted.
- Sometimes we are getting value greater than 1 & less than 0 (max/min are \therefore we need to squash this fn. $1 \& 0$)
- We use sigmoid activation fn.

Decision Boundary Logistic Regression

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$h_0(x) = \theta^T x$$

$$\text{let } z = \theta_0 + \theta_1 x$$

$$h_0(x) = g(z)$$

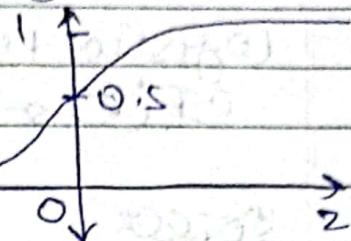
$$h_0(x) = \frac{1}{1+e^{-z}} \rightarrow \text{sigmoid or}$$

$$h_0(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}} \rightarrow \text{logistic fn}$$

$g(z) \geq 0.5 \quad \{$
when $z \geq 0 \}$

$g(z) < 0.5 \quad \{$
when $z < 0 \}$

$g(z)$



• Training set:

$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$y \leftarrow g_0, 1 \} \rightarrow (20/p)$$

$$h_0(z) = \frac{1}{1 + e^{-z}}, \quad z = \theta_0 + \theta_1 x$$

change parameter θ_1 .

C.F of linear Reg.

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^i) - y^i)^2$$

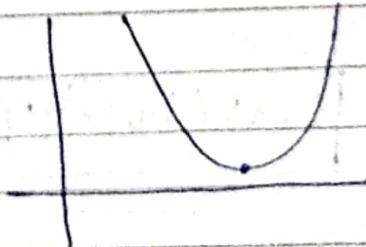
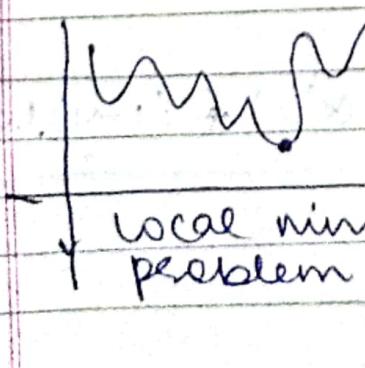
Logistic Regression,

$$C.F = \frac{1}{2} \sum_i (h_0(x^i) - y^i)^2 \quad] \text{we cannot use this as C.F.}$$

~~$J(\theta_0 + \theta_1 x^i - y^i)^2$~~

non-convex fn

convex fn



Logistic Regression CF.

~~Cost Function~~

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta_0 - \theta_1 x}}$$

Cost

$$J(\theta_0) = \begin{cases} -\log(h_{\theta}(x)), & y=1 \\ -\log(1-h_{\theta}(x)), & y=0 \end{cases}$$

↑ gives global min.

Cost = 0

$$\text{if } y=1, h_{\theta}(x)=1$$

Cost

if $y=0$

$$J(\theta_0)$$

0

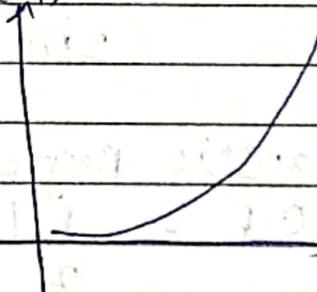
0

$h_{\theta}(x)$

(classification)

'if $y=0$

$$J(\theta_0)$$



$$\text{Cost}(h_{\theta}(x^i), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$$\text{Cost}(h_{\theta}(x^i), y) = -y \cdot \log(h_{\theta}(x)) + (1-y) \cdot \log(1-h_{\theta}(x))$$

$$J(\theta_0) = -\frac{1}{2m} \sum_{i=1}^m y_i \log(h_0(x_i)) + (1-y_i) \log(1-h_0(x_i))$$

Cost Fⁿ.

Convergence $\Delta \theta_0^m$

repeat until

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

?

Performance Matrix of Classification Problem

Actual Pred | Actual

X ₁	X ₂	Y ₁ Y ₂	{Pred}	1	0	
-	-	0 1		1	3	2
-	-	1 0		0	1	1
-	-	0 1				
-	-	1 0				

CONFUSION MATRIX

-	-	1 1	TP	FP	
-	-	1 1	FN	TN	
-	-	0 1			
-	-	1 0			

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} = 4/7 = 57\%.$$

① Precision $\rightarrow TP / (TP + FP)$

② Recall $\rightarrow TP / (TP + FN)$

③ F-score $\rightarrow 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$

ex: spam classification : Precision

person has cancer or not : Recall

stock market will crash : depends on prediction

↳ f-score

$$F\text{ score} = \frac{(1+\beta^2) \underset{\substack{\text{Precision} \\ \text{Recall}}}{\cancel{\text{Precision}}} \times \text{Recall}}{\beta^2 \times (\text{Precision} + \text{Recall})}$$

F1 score if $\beta=1$, $\frac{2(\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$

Harmonic Mean

$$\frac{2 \cancel{P.R}}{\cancel{P.R} + \cancel{R.P}}$$

F. S score $\beta=0.5$ $\text{FP} \geq \text{FN} \rightarrow \text{less important}$
 L more important

$$\frac{(1+\beta^2) P.R}{(\beta^2)(P+R)}$$

F2 score $\beta=2$. $\text{FN} > \text{FP} \rightarrow \text{less imp.}$
 L more imp

DAY-3 KNN & NAIVE BAYES ALGOM

Naive Bayes Inference (Classification)

↳ Bayes Th.

$$P(B|A) = \frac{P(B) \times P(A|B)}{P(A)} \rightarrow \text{Bayes' Theorem}$$

$\xrightarrow{\text{independent}}$ $\xleftarrow{\text{independent}}$

 $x_1 \ x_2 \ x_3 \ x_4 \dots x_n \ y$
~~P(Y|X)~~

$$P(Y|x_1, x_2, \dots, x_n) = \frac{P(Y) \times P(x_1, x_2, \dots, x_n|y)}{P(x_1, x_2, \dots, x_n)}$$

$$= \frac{P(Y) \times P(x_1|y) \times P(x_2|y) \times \dots \times P(x_n|y)}{P(x_1) \times P(x_2) \times \dots \times P(x_n)}$$

$y \begin{cases} \text{yes} \\ \text{no} \end{cases}$

$$P(\text{yes}|x_i) = \frac{P(\text{yes}) \times P(x_i|\text{yes}) \times \dots \times P(x_n|\text{yes})}{P(x_1) \times \dots \times P(x_n)} \quad \# \text{fixed}$$

$$P(\text{no}|x_i) = \frac{P(\text{no}) \times P(x_i|\text{no}) \times \dots \times P(x_n|\text{no})}{P(x_1) \times \dots \times P(x_n)} \quad \# \text{fixed}$$

\therefore we ignore denominator

$$P < 0.5 = 0$$

$$P(\text{yes}|x_i) = 0.13$$

$$P \geq 0.5 = 1$$

$$P(\text{no}|x_i) = 0.05$$

$$\text{Normalization} = \frac{0.13}{0.13 + 0.05} = 0.72 = 72\%$$

$$P(\text{no}) = 1 - 0.72 = 0.28 = 28\%$$

Dataset = $\{(\text{outlook}, \text{temp}, \text{humidity}, \text{wind}, \text{play})\}$

Outlook $P(\text{Sunny} | \text{yes})$

	Yes	No	$P(Y)$	$P(N)$
Sunny	2	3	2/9	3/9
Overcast	4	0	4/9	0/9
Rain	3	2	3/9	2/9
Total	9	5	14	14

Temperature

	$P(Y)$	$P(N)$
Hot	2	2/9
Mild	4	4/9
Cold	3	3/9
Total	9	5

Play

$$\text{Yes} - 9 \quad P(Y) = 9/14$$

$$\text{No} - 5 \quad P(N) = 5/14$$

Testcase $\rightarrow (\text{Sunny}, \text{Hot})$

$$P(\text{Yes/Sunny, Hot}) = P(\text{Yes}) * P(\text{Sunny}/\text{Yes}) * P(\text{Hot})$$

$$= \frac{9}{14} * \frac{2}{9} * \frac{2}{9}$$

$$= \frac{2}{63} = 0.031$$

$$P(\text{no/Sunny, hot}) = \frac{9}{14} * \frac{3}{9} * \frac{2}{5} = \frac{3}{35} = 0.085$$

$$P(Y_{CS}/x_i) = 0.031 + 1.78 = 27\%$$

$$P(NO/x_i) = 0.085 = \frac{0.085}{0.085+0.031} = 73\%$$

$\rightarrow (\text{sunny}, \text{hot}) \rightarrow \text{NO}$
 ↳ pred.

$\rightarrow (\text{overcast}, \text{med})$

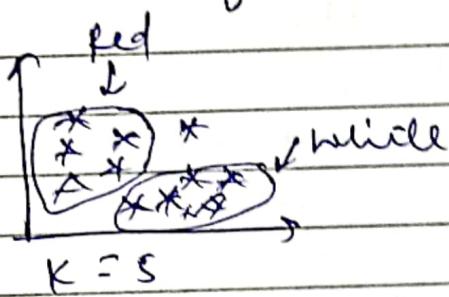
$$P(Y/x_i) = \frac{2.4}{9} \times \frac{9}{7.4} \times \frac{4}{9} = \frac{8}{63}$$

$$P(N/x_i) = 0$$

KNN Alg^m

- ↳ classification
- ↳ regression

(i) ↳ classification



$$\text{max no red} = 3 \\ \text{white} = 2 \}$$

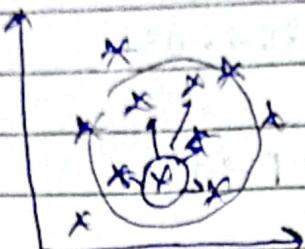
we calculate dist using euclidean dist & manhattan dist.

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

$$\text{Manhattan} =$$

$$\sqrt{|(x_2 - x_1) + (y_2 - y_1)|}$$

(ii) Regression



$$k = 5$$

↳ hyperparameter
Avg of all pts is o/p.

KNN Worst cases

(i) Outliers

(ii) Imbalanced dataset

Difficulties in handling imbalanced datasets
Outliers and noise

Difficulties in handling imbalanced datasets
Outliers and noise

DAY-4 DECISION TREE & ENSEMBLE ALGO'S

1) Decision Tree

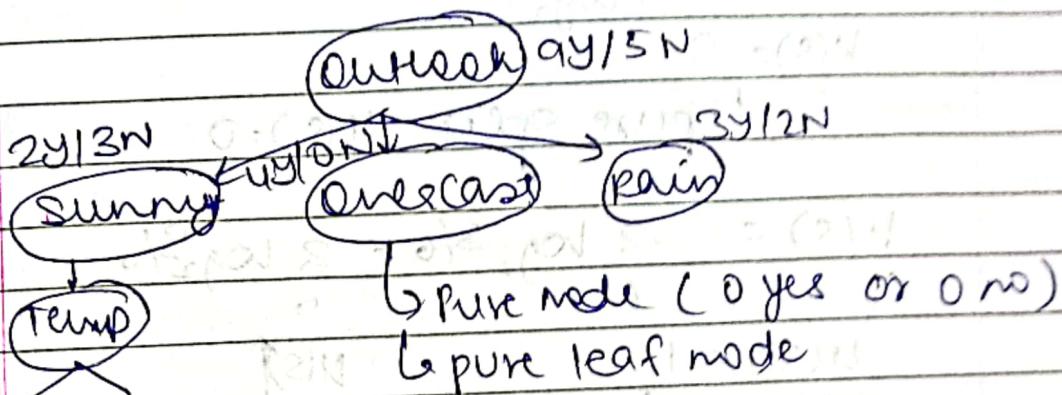
- ↳ Regression
- ↳ Classification

we use nested if-else conditions

(i) Classification

Dataset.

day outlook temp humidity wind play



~~pure option~~
we will split till we get a leaf node.

① Purity:

↳ entropy

↳ Gini coefficient / impurity

② How the features are selected?

↳ Information gain

③ Entropy

$$H(p) = -p \log_2 p - (1-p) \log_2 (1-p)$$

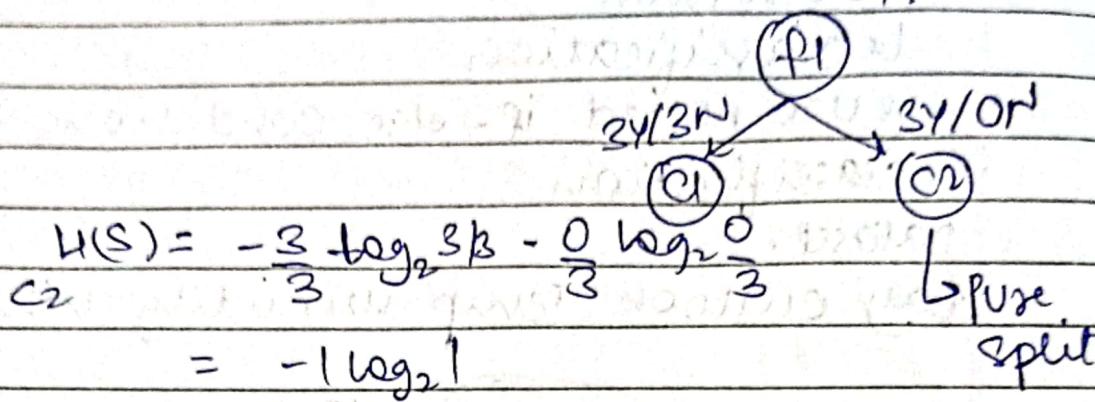
④ Gini Impurity

① Entropy

$$H(S) = -P_1 \log_2 P_1 - P_0 \log_2 P_0$$

$P_i \rightarrow$ Prob. of Yes

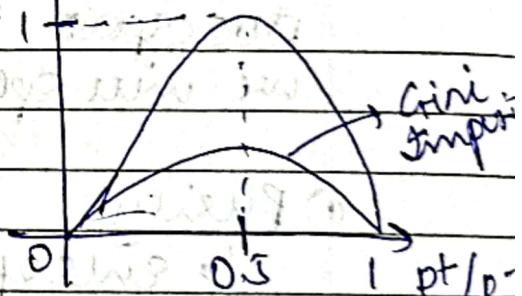
64/3N



$$H(S) = -\frac{2}{3} \log_2 \frac{2}{6} - \frac{1}{3} \log_2 \frac{1}{6}$$

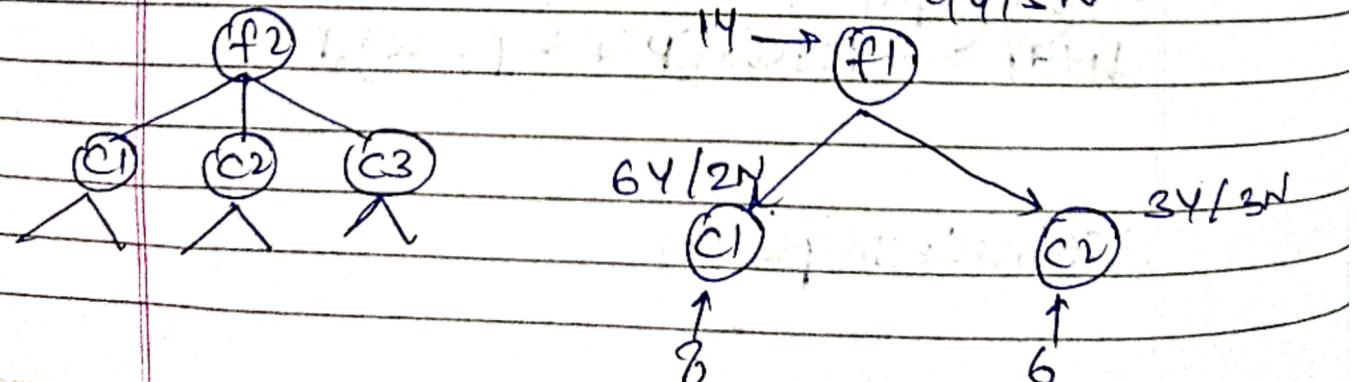
$$H(S) = 1.392 \text{ bits}$$

$$0 \leq H(S) \leq 1$$



Purity test is done with entropy

② Which feature to split



Information Gain

$$\text{Gain}(S, f_1) = H(S) - \sum_{\text{val}} |S_{\text{val}}| \cdot H(S_{\text{val}})$$

$$H(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} \approx 0.94$$

(approx)

$$H(S_{v1}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \approx 0.81$$

$$H(S_{v2}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \approx 1$$

$$\text{Gain}(S, f_1) = 0.94 - \left[\frac{8}{14} * 0.81 + \frac{6}{14} * 1 \right]$$

$$\text{Gain}(S, f_1) = 0.049$$

$$\text{Gain}(S, f_2) = 0.051$$

$$G(S, f_2) \geq G(S, f_1)$$

\therefore we will start the feature split with f_2 .

② Gini Impurity

$$G_I = 1 - \sum_{i=1}^n (p_i)^2$$

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

$$= 1 - [(V_1)^2 + (V_2)^2]$$

$$= 0.5$$

$$n=2 \quad O/p \rightarrow \text{Yes}$$

$$2/2N$$

6. snap

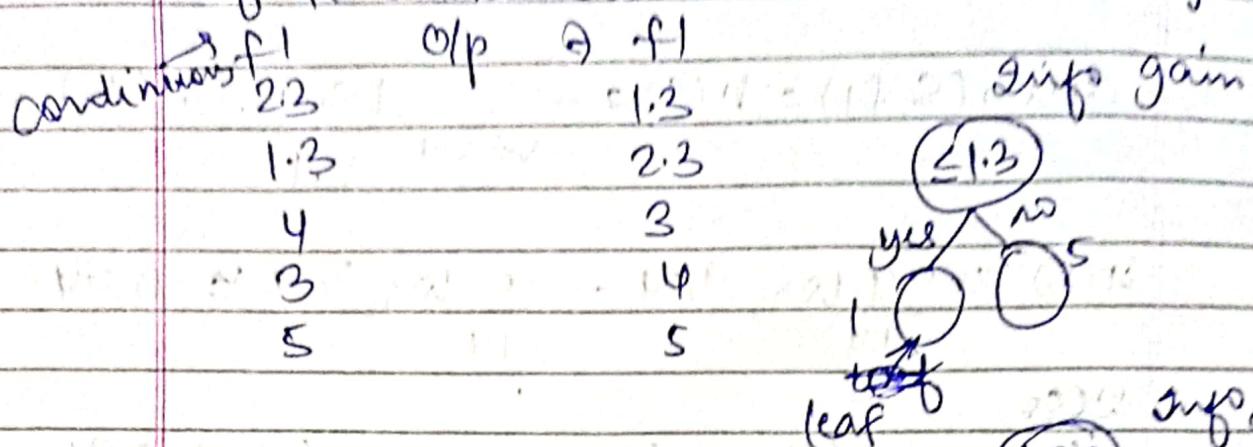
Impure split

$$0 \leq G_I \leq 0.5$$

x. Gini Impurity is faster than entropy.

No. of features is greater \rightarrow gini less \rightarrow entropy

If f_1 is numerical rather than object.



It will cal

Information gain

for all the records and the best will be selected.

Decision Tree regressor

$f_1 \rightarrow$ continuous

↳ mean MSE or MAE

$$\text{mean} = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

↳ mean \hat{y}_i

In regression instead of using entropy or gini impurity we use MSE or MAE.

Hyperparameters

Decision → Overfitting

↳ Post pruning

↳ Pre pruning

(max-depth, max-leaf)

↳ GridSearchCV

Random Forest

Decision Tree

Random Forest

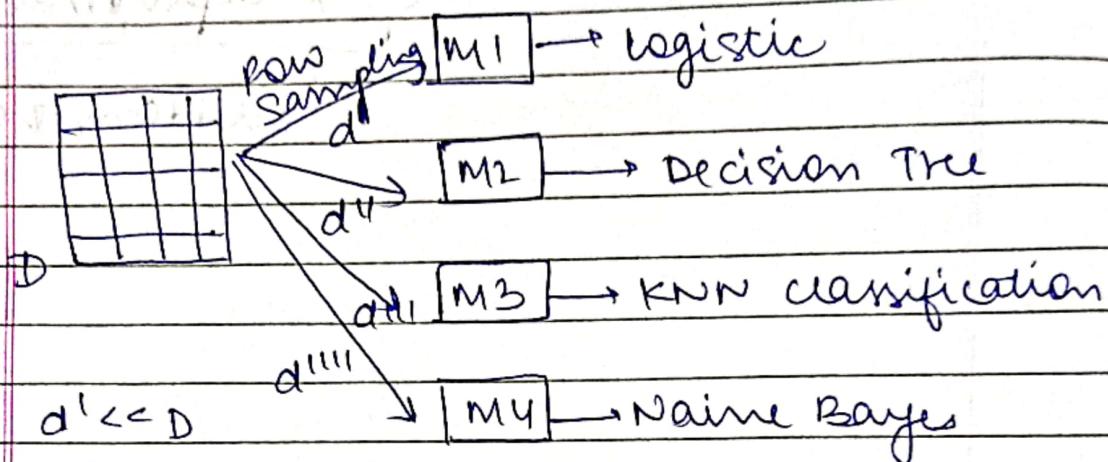
Decision Tree

DAY:5 AdaBoost, Random Forest, XGBoost

Ensemble Techniques

↳ Multiple algos to solve a problem.

(i) Bagging - Classification



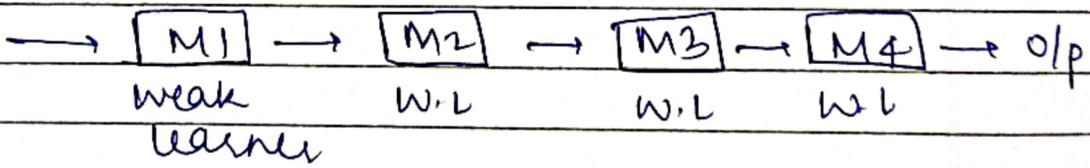
- The records may repeat.
- There will be 100 to 200 models.

- Regression

- Average/mean will be taken.

\Rightarrow

(ii) Boosting



$\frac{1}{4}$ strong learner
 ↳ pred is bad

- It's a sequential set of all the models combined together and these models are weak learners.

↳ If we combine multiple W.L then we have strong learner.

↳ If we combine multiple W.L then we have strong learner.

when we combine them they give us strong learners.

BAGGING

- 1. Random Forest classifier
- 2. Random Forest Regressor

BOOSTING

- 1. AdaBoost
- 2. Gradient
- 3. Xgboost

1. Random Forest classifier

Q. What is the main disadvantage of decision tree?

↳ Overfitting

\therefore Low bias & high variance.

- Normalization is not required in decision tree.
- Random forest is not affected by one outliers

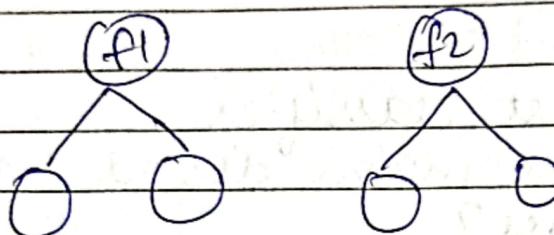
Adaboost

f1	f2	f3	f4	O/P	weight
-	-	-	-	Yes	Y ₁
-	-	-	-	No	Y ₂
-	-	-	-	-	Y ₃

Σ Y

Overall sum of w = 1

- we create decision trees by selecting the features with the help of information gain or gini impurity.
- In random forest we create the decision trees upto only one levels, known as STUMPS.



{weak learners}

~~Step I: If total error = 0 then stop~~

Step I: If only one record was predicted wrong.

$$\therefore \text{Total error} = Y_T$$

Step II. Performance of stump

$$= \frac{1}{2} \log_e \left(\frac{1 - T.E}{T.E} \right) \rightarrow \text{Total error.}$$

$$= \frac{1}{2} \log_e \left(\frac{1 - Y_T}{Y_T} \right)$$

$$\approx 0.895$$

Step III. update all the weights \uparrow performance of new sample wt. \uparrow Stump

For correct records: $wt \times e^{-P_S}$

$$= \frac{1}{7} \times e^{-0.895} \approx 0.05$$

For incorrect records: $wt \times e^{P_S}$

$$= \frac{1}{7} \times e^{0.895} \approx 0.349$$

- The wt of incorrect records are higher than the correct records b/c. if they have lower weights then they will train more by passing them through more weak learners.

$$\sum \text{New wt} \neq 1$$

but we know the wt is always equal,

: we find the normalized wt

divide all the wt. by the ~~sum of new. wt.~~
sum of new. wt.

new wt: 0.05, 0.05, 0.05, 0.349, 0.05, 0.05, 0.05

normalized wt: $\frac{0.05}{0.649}, \frac{0.05}{0.649}, \frac{0.349}{0.649}, \dots$

: $0.07, 0.07, 0.07, 0.537, \dots, 0.07$

now ~~too~~ buckets will be created

$[0 - 0.07], [0.07 - 0.14], [0.14 - 0.21],$
 $\underline{[0.21 - 0.747]}, [0.747 - 0.754], \dots$

~~6 strong tag~~

- In classification problem we take the majority ans.
- In regression we take the mean / avg of all the pred / o/p.

- Black Box Model: Here we are not able to see how all the weights are updated. The maths behind the algm is complex.
ex: Random Forest, ANN
- White Box Model: The maths is simple we can visualize how all the weights are updated.
ex: Linear Regression, Decision Tree

Explain what is meant by white box model
and black box model.

Ans: White box model is one in which we can understand the working of the model.

Explain what is meant by black box model.

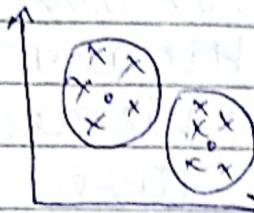
Ans: In black box model we cannot understand the working of the model.

DAY-6 Kmeans, Hierarchical clustering, Silhouette score, DBSCAN clustering

#. Kmeans clustering

L_k = centroids

for $k=2$



Step I: we try with diff. k values

for $k=1, 2, \dots, n$

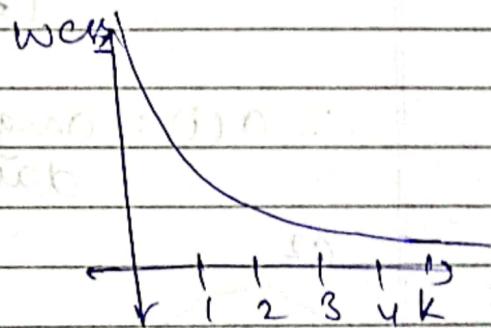
Step II: initialize 'k' no. of centroids

Step III: compute the avg. w/ update centroids.

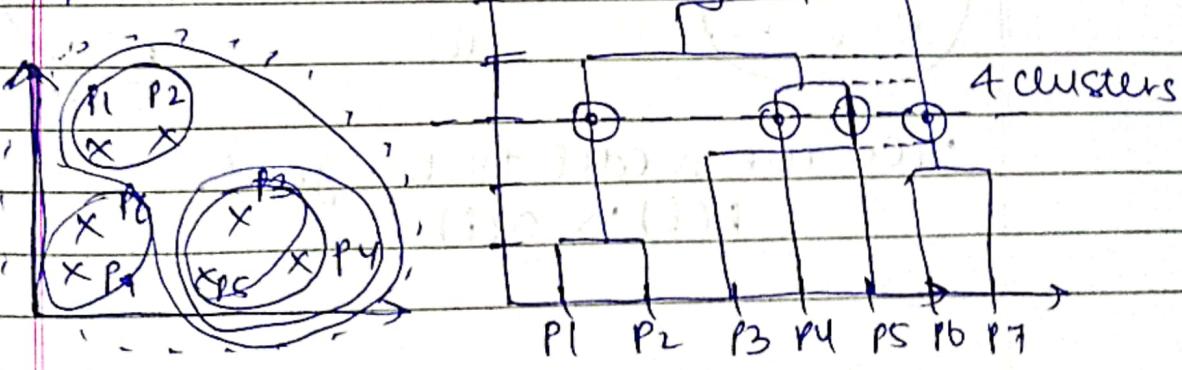
* Elbow Method (K value)

WCSS - within cluster sum of square

for $i=1-10$



Hierarchical clustering \rightarrow Dendrogram



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

Q. What time is taken by k-means or Hierarchical clustering?

→ Hierarchical clustering will take more time.

- If d_s is small → ~~Hierarchical~~ Hierarchical
- If d_s is large → k-Means

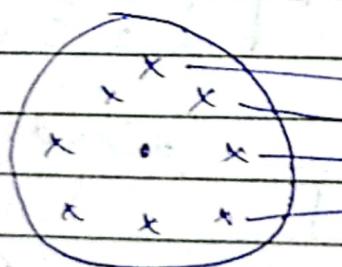
Validation of Clustering Model
↳ Silhouette score

$$a(i) = \frac{1}{(C_i - 1)} \sum_{j=0}^{C_i - 1} d(i, j)$$

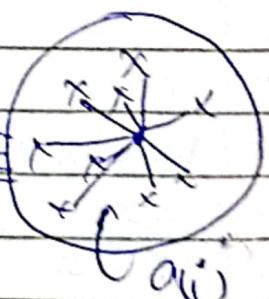
$\underbrace{(C_i - 1)}_{\text{dist bet}} \quad \underbrace{\sum_{j=0}^{C_i - 1} d(i, j)}_{\text{dist bet between centroid & data pts}}$

$\therefore a(i)$ = avg of dist bet centroid & data points.

c2



b(i)



c1

a(i)

- For the model to be good
 $b(i) \gg a(i)$

$$b(i) = \min_j \frac{1}{C_j} \sum d(i, j)$$

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

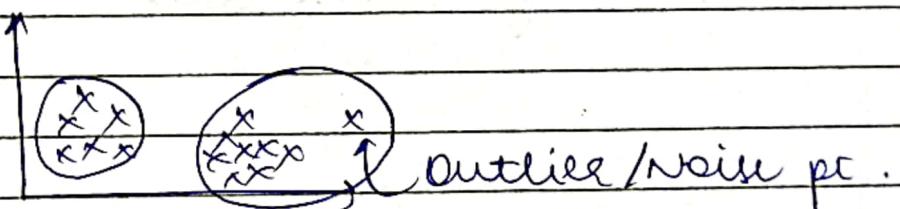
$$-1 \leq s(i) \leq 1$$

the better model

DBScan clustering

(Density Based spatial clustering of Application with noise)

- We leave the outliers separately and these outliers are known as noise pts.



- Epsilon radius of the circle centered at point x created by ~~other~~ ~~not~~ ~~the~~ ~~other~~ points in the neighbourhood around a pt x .
 - minpts : The density threshold.
If a neighbourhood includes at least minpts. it will be considered as a dense region. Alternatively, a point

will be considered as dense if there are atleast the value of minpt in its ϵ -neighbourhood. These dense pts are called core points.

- Border Point: has ϵ -neighbourhood that contains less than minpt but it belongs to the ϵ -neighbourhood of another core pt.

Bias is a phenomenon that skews the results of an algm in favor or against an idea (training ds.).

Variance refers to the changes in the model when using different proportion of training or test data.

DAY-7 XGBoost, SVM, SUR

→ Extreme Gradient Boosting

Xgboost classifier:

	Salary	Credit	Approval	Residual
≤ 50	B	0	-0.5	
≤ 50	A	1	0.5	
≤ 50	A	1	0.5	
> 50	B	0	-0.5	
> 50	A	1	0.5	
$> 50K$	BN	1	0.5	

Base Model $\rightarrow \text{Pr} = 0.5$

\hookrightarrow probability of approval

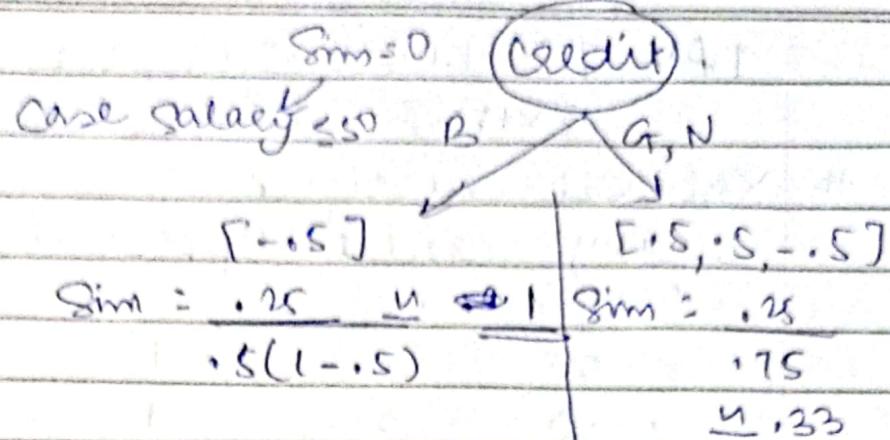
Step I: Create a binary decision tree using the features.

Step II: Cal Similarity wt $= \frac{\sum (\text{Residual})^2}{\sum (\text{Pr}(1-\text{Pr}) + \lambda)}$

Step III: Calculate Information gain

$$\begin{array}{c}
 [-0.5, +0.5, +0.5, -0.5, +0.5, +0.5, -0.5] \\
 \text{Salary} \hookrightarrow 0.14 \\
 \text{Sp. } \leq 50 \quad \checkmark \quad \geq 50 \\
 \frac{(-0.5 + 0.5 + 0.5 - 0.5)^2}{4(0.5)(1-0.5)} \quad \frac{(-0.5 + 0.5 + 0.5)^2}{3(0.5)(1-0.5)} \\
 \approx 0.222222 \approx 0.33
 \end{array}$$

$$\text{InfoM gain} = 0 + 0.33 - 0.14 = 0.19$$



$$\text{Inform gain} = 1 + 0.33 - 0 = 1.33$$

* $\sigma [0 + \alpha(\text{sim})]$
 ↳ learning rate
 ↳ sigmoid fn.

Q. $\sum_{i=0}^n \sigma_i$

• xgboost → black box Model

xgboost Regression

Expn	Gap	Salary	RPs
2	4	40K	-11K
2.5	4	48K	-9K
3	N/A	52K	1
4	N	60K	9
4.5	4	62K	11

Base Model = avg o/p = 51K

(-11, -9, 1, 11)

Eppn

$\text{sim} \leq 46$

≤ 2

> 2

(-11)

(-9, 1, 9, 11)

$\text{Sim} = \frac{1}{2} \ln 605$

$\text{sim} \leq 28.8$

$\text{Sim}_{\text{WT}} = \sum (\text{residual})^2$

No. of residual + 1

↳ Hyperparameter

Info^m gain = 60.5 + 28.8 - 46 ≤ 42.3

(-11, -9, -10) ≤ 2.5 > 2.5
 Eppn $\rightarrow \text{sim} \leq 46$
 [-11, -9] [1, 9, 11]

$\text{sim} \leq 133.33$ $\text{sim} \leq 110.25$

Info^m gain = 133.3 + 110.25 - 46

$$\mathcal{O}(p=51 + \alpha_1(-10) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n))$$

Hard/Soft \rightarrow in case of overlapping
marginal plane

SVM.

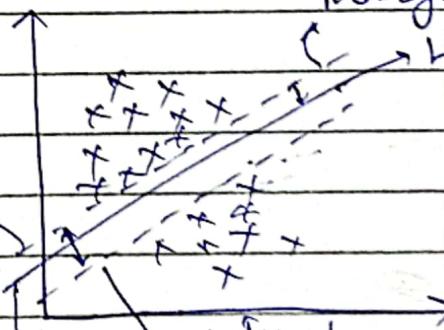
$$w^T x + b = -1$$

$$w^T x_1 + b = 1$$

$$(-) w^T x_2 + b = -1$$

$$\frac{w^T(x_1 - x_2)}{|w|} = 2$$

$$\frac{|w|}{|w|}$$



$$\therefore y = mx + c$$

we have to maximise $2/|w|$.

such that,

$$y_i \begin{cases} \geq +\epsilon, & w^T x + b \geq 1 \\ = 1, & w^T x + b = 1 \\ \leq -\epsilon, & w^T x + b \leq -1 \end{cases}$$

or

$$\underline{y_i + (w^T x_i + b) \geq 1}$$

for correct points

$$\max_{(w,b)} \frac{2}{\|w\|} \Rightarrow \min_{(w,b)} \frac{\|w\|}{2}$$

$$\min_{(w,b)} \frac{\|w\|}{2} + C \sum_{i=1}^n \eta_i$$

to sum of dist. of
the wrong pts.

SVR

is similar to SVM with some minor changes in the part

$$\underline{C \sum_{i=1}^n \eta_i}$$

SF