

Actual Predictions

1	1	✓
0	0	✓
1	0	✓
1	1	✓
0	1	
0	1	
0	0	✓
1	1	✓
1	0	
0	1	

Actual Accuracy = 50%

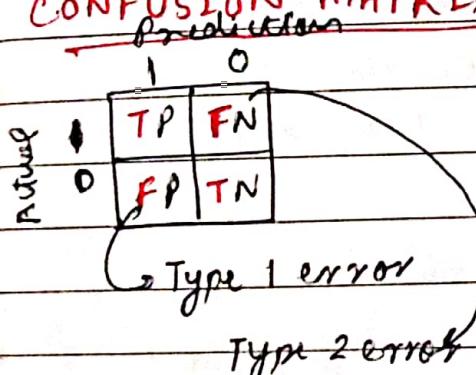
① Accuracy depends on the content of the problem.

② Accuracy is not ^{good} with imbalance data.

Spam message is say 0.1% of total data & it is predicting for correct predictions for all;
Accuracy = 99.9% (No use)

③ Accuracy considers $1 \rightarrow 0$ & $0 \rightarrow 1$ predictions as same.

CONFUSION MATRIX



		Predictions	
		1	0
Actual	1	3	2
	0	3	2

		Predictions	
		Spam	Not Spam
Actual	Spam	100	170
	Not Spam	100	700

Accuracy (M_1) = ~~800~~

Accuracy (M_2) = 80%

		Predictions	
		Spam	Not Spam
Actual	Spam	100	190
	Not Spam	100	700

But we will take M_2

Reason: ~~if~~ non-spam ~~is~~ spam ~~it's~~ problem \nrightarrow
(Imp. document can be lost)

$$(FP)_{M_1} > (FP)_{M_2} \Rightarrow (\text{type-1 error})_{M_1} > (\text{type-1 error})_{M_2}$$

Precision:- What proportion of predicted positive is truly Positive,

$$\text{Precision} = \frac{TP}{TP + FP}$$

We can also calculate precision for any class (0 or 1)

[Precision for each class can be calculated]

		Predictions		Predictions	
		Detected Cancer	No Cancer	Detected Cancer	No Cancer
Actual Cancer	Detected Cancer	1000	(200) FN	Detected Cancer	1000
	No Cancer	800	8000	No Cancer	500

M_1 M_2

Accuracy (M_1) = Accuracy (M_2) type-2-error

$$\Delta \text{ Recall} = \frac{1000}{(M_1) 1000 + 200} \rightarrow \text{Recall}(M_2) = \frac{1000}{1000 + 500}$$

Better to take M_1

दिए दिए सभी डेटा को नहीं मार्गित करते हैं।

Recall :- Proportion of actual positive is correctly classified.

Recall = $\frac{TP}{TP+FN}$	⊗ ⊗ ⊗
-----------------------------	----------

We can calculate recall of all individual cases.

Input image is cat or dog. precision or

Then problem is whether to consider recall.

Use F1 Score.

F1 Score = harmonic mean of precision & recall. → brings F1 to lower value among two

F1 Score = $\frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$	⊗ ⊗ ⊗
--	----------

	M_1	M_2
Precision	80	60
Recall	80	100

$$\text{Figure} \Rightarrow .80 \quad .75$$

$$F1(M_1) \quad F1\text{Score}(M_1) > F1\text{Score}(M_2)$$

i.e., M_1 is better than M_2

MULTICLASS CLASSIFICATION :-

predicted

	Dog	Cat	Rabbit	
Dog	25	5	10	= 40
Actual Cat	0	30	4	= 34
Rabbit	4	10	20	= 34

$$\text{Accuracy} = \frac{25 + 30 + 20}{40 + 34 + 34} \times 100 = 69.4\%$$

$$\text{precision (Dog)} = \frac{25}{25+0+4} \times 100 = 86.20\%$$

$$\text{precision (Cat)} =$$

$$\text{precision (Rabbit)} =$$

Macroprecision :-

Macroprecision : If dataset is balanced then calculate average of all ~~total~~ precision.

$$\text{Precision (Model)} = \frac{P(\text{dog}) + P(\text{cat}) + P(\text{rabbit})}{3}$$

If not balanced then use weighted precision (SK way)

$$\text{Precision (Model)} = \frac{40}{108} \times P(\text{dog}) + \frac{34}{108} P(\text{cat}) + \frac{34}{108} P(\text{Rabb})$$

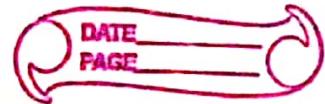
$$\text{Recall (Dog)} = \frac{25}{25+5+10} = 62.5\%$$

Same as precision
{ average }
{ & weighted }

$$F_{1\text{ score}} (\text{Dog}) = \frac{2 * \text{precision}(\text{Dog}) * \text{Recall}(\text{Dog})}{\text{precision}(\text{Dog}) + \text{Recall}(\text{Dog})}$$

same as
precision
{ average &
weighted }

9.02.2024
Monday



NAIVE BAYE'S CLASSIFIER

- ① Supervised learning algorithm used for classification.
- ② It is a part of generative learning algorithm.

Learning Algorithm

Generation

Two is
one's
assumption for. ② Naive Baye's

Learn the features of class
ex. Car & Dog is there

To which is more similar to?

discriminative learning Algo.

① Logistic Regression ③ Neural
② Support Vector Machines N/W

Things that differentiate two
classes.

This classifier works on the principle of conditional probability.

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

Prob. of A given B already happened

Q) Suppose a coin is tossed 2 times. What is the prob. that
both the toss is heads given that atleast one of the toss is head.
Sol. $S = \{HH, HT, TH, TT\}$

$$B = \{HH, HT, TH\}$$

$$A = \{HH\}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{1/4}{3/4} = \frac{1}{3}$$

Bayes Theorem :-

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \Rightarrow P(A \cap B) = P(A|B) \cdot P(B) \rightarrow \textcircled{I}$$

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

$$\Rightarrow P(A \cap B) = P(B|A) \cdot P(A) \rightarrow \textcircled{II}$$

From \textcircled{I} & \textcircled{II} :

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

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

↳ Summary:
Every feature is
independent.

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

↑ likelihood
↑ prior

$P(B)$ → Marginal prob. probability

Likelihood Prob. :- It is the relative probability of the observations occurring for each other. "Posterior probability."

Marginal Prob. :- Prob. of observation

Prob. of the event after the evidence has been

Name Classification Algo.

Suppose we have a dataset of n features $(x_1, x_2, x_3, \dots, x_n)$:
 $\rightarrow y$ is o/p. variable.

$x_1, x_2, x_3, \dots, x_n$	y	<i>miss fit</i>
classroom	b_1 b_2 \vdots	Find prob.

$$\text{Yes } P(c_1|x) = P(y_{\text{yes}}|x) = \frac{P(x|c_1) \cdot P(c_1)}{P(x)}$$

$$\text{No } P(c_2|x) = P(y_{\text{no}}|x) = \frac{P(x|c_2) \cdot P(c_2)}{P(x)}$$

D^T does not play importance.

$$\text{Output Class} = \underset{c_k}{\operatorname{argmax}} \{ P(c_k) \oplus P(x|c_k) \}$$

$$= \underset{c_k}{\operatorname{argmax}} (P(c_k) \cdot P(x|c_k))$$

Org. \oplus different or \oplus values given

$$P(x_i | c_k) = P(x_1, x_2, \dots, x_n | c_k)$$

~~Naive Assumption~~

Naive Assumption

No pair of ~~same~~ features are dep.

$$P(x_1, x_2, \dots, x_n | c_k) : P(x_1 | c_k) * P(x_2 | c_k - \text{ & } \dots)$$

(All assumed)

& similarly we can do for x_n



$$P(x_i | c_k)$$

$$P(x_1, x_2, \dots, x_n | c_k) = P(x_1, x_2, \dots, x_n | c_k) / n$$

P.C. \rightarrow c_k

input class = argmax $\{P(c_k) \prod_{i=1}^n P(x_i | c_k)\}$

Play Tennis

		features				
Say	outlook	Temp	Humidity	wind	PlayTennis	
1	Sunny	Hot	High	weak	No	
2	Sunny	Hot	High	Strong	No	
3	Overcast	Hot	High	Weak	Yes	
4	Rain	Mild	High	Weak	yes	
5	Rain	Cool	Normal	Weak	Yes	
6	Rain	Cool	Normal	Strong	Yes No	
7	Overcast	Cool	Normal	Strong	No Yes	
8	Sunny	Mild	High	Weak	Yes No	
9	Sunny	Cool	Normal	Weak	No Yes	
10	Rain	Mild	Normal	Weak	Yes	

~~20.02.2024
Tuesday~~

$$\arg\max \left\{ \begin{array}{l} P(C_1 | \text{<end>}), \\ P(C_2 | \text{<end>}), \\ P(C_3 | \text{<end>}) \end{array} \right\}$$



11	Sunny	mild	Normal	Strong	Yes
12	overcast	mild	High	Strong	Yes
13	overcast	Hot	Normal	weak	Yes
14	Raining.	mild	High	Strong	No

Q) < sunny, ~~Hot~~ Hot, Normal, Strong > for a day.
What are we going to play tennis or not?

Will we go to play tennis or not?

Sol: Lookups tables; learning process } conditional probability

$$\text{Output Class} = \underset{c_k}{\operatorname{argmax}} \left\{ P(c_k) * \prod_{i=1}^n (x_i | c_k) \right\}$$

↳ class

~~veryman~~ $\{ P(\text{yes} \mid \{\text{sunny, hot, normal, strong}\})$
 $P(\text{no} \mid \{\text{sunny, hot, normal, strong}\}) \}$

Assuming that all the features are independent

$$P(\text{Yes}) = P(\text{Sunny} | \text{Yes}) * P(\text{Hot} | \text{Yes}) * P(\text{Normal} | \text{Yes}) * P(\text{Strong} | \text{Yes})$$

$$\rightarrow P(\text{No}) * \underbrace{P(\text{Sunny} | \text{No}) * P(\text{hot} | \text{No}) * P(\text{Normal})}_{\text{Prior probability}} * \underbrace{P(\text{Rainy} | \text{No}) * P(\text{Cold} | \text{No}) * P(\text{Unusual})}_{\text{Conditional probability}}.$$

$$P(\text{Yes}) = 9/14 \quad P(\text{No}) = 5/14$$

Likelihood Probability Tables — Lookup Tables.

Outlook

	Year	No
Sunny	2/g	3/5
Overcast	4/g	0/5
Rain	3/g	2/5

$$P(\text{outlook} = \text{sunny} \mid \gamma_{20}) = \frac{P(\text{sunny} \cap \gamma_{20})}{P(\gamma_{20})}$$

Temp

	Yes	No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

Wind

	Yes	No
Weak	6/9	2/5
Strong	3/9	3/5

Humidity

	Yes	No
High	3/9	4/5
Normal	6/9	1/5

$$\text{argmax} \left\{ \frac{9}{14} \times \frac{2}{3} \times \frac{1}{3} \times \frac{6}{9} \times \frac{3}{9} = \frac{36}{81} = 0.007 \right.$$

$$\left. \frac{5}{14} \times \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{3}{5} = 0.01 \quad \checkmark \right.$$

Probability for 'No' is higher $\therefore \underline{\text{output}} = \text{No}$.

Q) input <Overcast, Hot, Normal, Strong>

Sol: $P(\text{No} | \text{input}) = 0 \quad P(\text{Overcast} | \text{No}) = 0$
 But it is not logical.

Since dataset is small.

"Zero Probability Error"

Solution: Add atleast one hypothetical data entry

"Laplace Smoothing"

$$P(x_{ij}|c) = \frac{P(x_{ij}|c) + \alpha}{N + \alpha K}$$

$\alpha = 1$ \downarrow_{new}

$\alpha = \text{smoothing parameter}$
 (hyperparameter)

$N = \text{No. of test examples}$

$$P(\text{Overcast} | \text{No}) = \frac{0+1}{5+1 \times 3} = \frac{1}{8}$$

$\frac{1}{8}$ | No

with o/p = c

$K = \text{No. of possible distinct}$

values for feature j.

\hookrightarrow Overcast $\xrightarrow{3 \text{ w classes}}$
 $\xrightarrow{\text{Sunny Overcast, Rainy}}$

2) formula

\equiv

Let zero del change del 2)

21.02.2024
Wednesday



$$\begin{aligned}
 P(\text{Sunny} / \text{No}) &= \frac{3+1}{5+1 \times 3} = \frac{4}{8} \\
 P(\text{Rain} / \text{No}) &= 3/8 \\
 P(\text{Overcast} / \text{No}) &= \frac{1}{8}
 \end{aligned}
 \quad \left. \begin{array}{l} \text{optional to} \\ \text{change} \\ \text{for some.} \end{array} \right\} \sum_{\text{Epi}} = 1$$

$$\therefore P(\text{No}) = 8/17 \quad P(\text{Yes}) = 9/17 \quad (\text{Not compulsory for two change})$$

$$\left. \begin{array}{l} P(\text{No}) + P(\text{O}/\text{No}) + P(\text{H}/\text{No}) + P(\text{N}/\text{No}) + P(\text{Str.}/\text{No}) \\ \frac{8}{17} \times \frac{1}{8} \times \frac{2}{5} \times \frac{1}{5} + \frac{3}{5} \end{array} \right\} = 2.823 \times 10^{-3} = 0.0028$$

avg risk

$$\left. \begin{array}{l} P(\text{Yes}) + P(\text{O}/\text{Yes}) + P(\text{H}/\text{Yes}) + P(\text{N}/\text{Yes}) + P(\text{Str.}/\text{Yes}) \\ \frac{9}{17} \times \frac{4}{9} \times \frac{2}{9} \times \frac{6}{9} \times \frac{3}{8} \end{array} \right\} = 0.011$$

α की higher value पर "underfit" हो जायेगा।
 α की lower value पर "overfit" हो जायेगा।

Advantages:-

(1) Simple to implement & use.

(2)

Disadvantage:-

(1) zero probability error.

(2) All situations can't be modelled.

DECISION TREE

Non linear

① $\begin{cases} \hookrightarrow \text{Classification} \\ \hookrightarrow \text{Regression} \end{cases}$

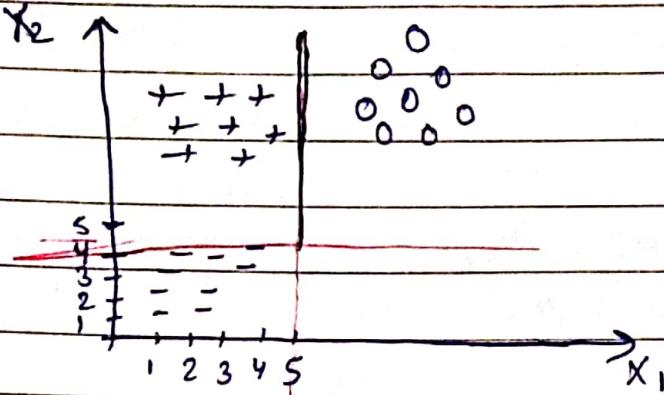
② Non linear model

③ Dividing feature space on the basis of feature

"Decision trees works by dividing the data in regions based on "if-then" type of questions."

④ Graphically by asking many "if then" questions on decision tree

can divide up the feature space using little segments of vertical & horizontal lines.



$$\text{if } (x_2 \leq 4)$$

print (" - Class ")

else if ($x_1 \leq 5$)

print (" + Class ")

Q) Gender Occupation Suggestion

F Student PUBG

F Programmer Github

M programmer Whatsapp

F programmer Github

M student PUBG

M student PUBG

Sol: if (occupation = 'Student')

print ("PUBG")

else

if (Gender = 'F')

print ("Github")

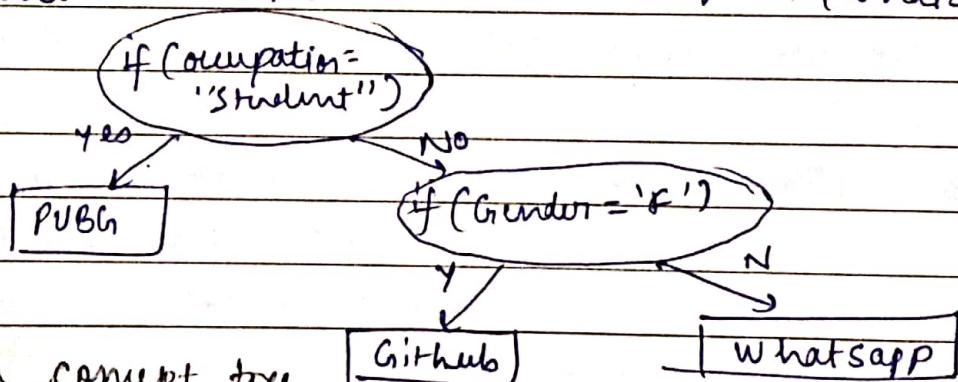
else

print ("Whatsapp")

Sol: Control

Flow

Diagram



Decision Tree is a concept tree

Github

Whatsapp

which summarizes the information contained in the training dataset in the form of a tree structure.

Structure of a Decision Tree

i) Root node

→ Top most node of tree is called root node

ii) Internal node / Decision node → Decision Node represents a test condition of an input feature & off of this test condition are branches come out from the decision node.

iii) Branches

iv) Leaf node

\Rightarrow Recursive greedy algorithm

is Decision Tree Algo.

\Rightarrow Decision tree is representing the complete hypothesis test.

distinct

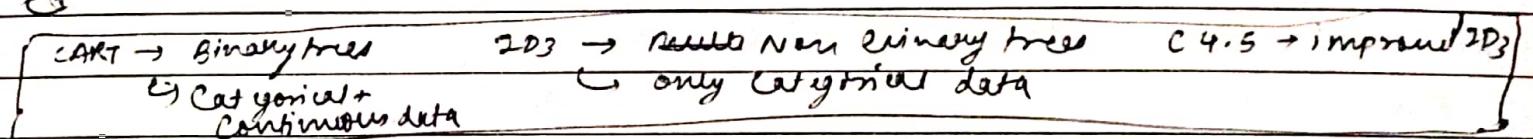
Max. branches from node = Max. no. of feature values a feature can have.

Decision node can have maximum out degree (branches) equal to possible values of features used to test in the node.

Termination Condition :-

Leaf node represents the outcome of the decision path.
Labels of the leaf nodes are the different classes a data point can belong to.

Q



Training process

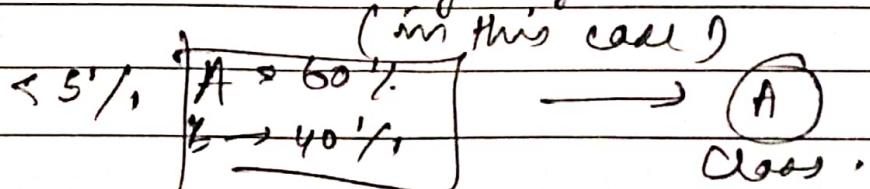
- ① The tree is constructed in a top-down fashion, & it starts from the root node.
- ② At the root we have entire dataset and we need to find the test split attribute.
- ③ After split each child represents splitted dataset based on values of the features.
- ④ This process is recursive and continued until we end up in the last level of the tree or finding a leaf tree.

\hookrightarrow we can decide previously only that we need only 20 levels only.

3 major stopping criteria :-

- ① Leaf Node
- ② Node already decided (say 20 levels i.e. total)
- ③ If data at any level less than say 5%, data (say) then stop division \rightarrow majority wins

Else we will get overfitted tree.



Decision-tree Induction Algorithms :-

① ID3 :-

This algorithm uses "Entropy" to measure the purity / impurity of the dataset at a node. (leaf node)

$$\text{Entropy} = - \sum_k P(c_k) \cdot \log_2(P(c_k))$$

For 2 classes - Entropy = $-P(c_1) \cdot \log_2(P(c_1)) - P(c_2) \cdot \log_2(P(c_2))$

Single class \rightarrow Entropy

For a 2 class problem, the min. entropy is 0 and max is 1
~~- single 100%, 0%~~, 50%, 50%

For more than 2 classes the min. entropy is 0 but max can be greater than 1. (Refer to first log base 2, 3, 4, 5)

| Information Gain = Entropy before split - Entropy after split

| How much entropy is reduced on splitting on basis of this.

Entropy after split = weighted sum of entropies of children

$$\text{Weight of child} = \frac{\text{No. of examples in child node}}{\text{total No. of examples}}$$

Q) Color Diameter Fruit Sol:

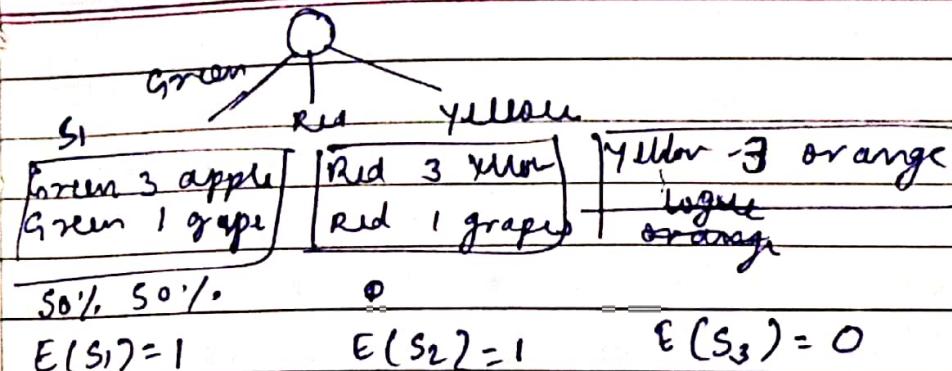
Green	3	apple
Red	3	apple
Red	1	grapes
Yellow	3	orange
Green	1	grapes

$$E(S) = - \sum_k P_k \log_2(P(c_k))$$

$$= \sum_{i=1}^n P_i(c_i) \cdot \log_2(1 - P_i(c_i))$$

$$= \left\{ \frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \frac{2}{5} \log_2 \left(\frac{3}{5} \right) + \frac{1}{5} \log_2 \left(\frac{1}{5} \right) \right\}$$

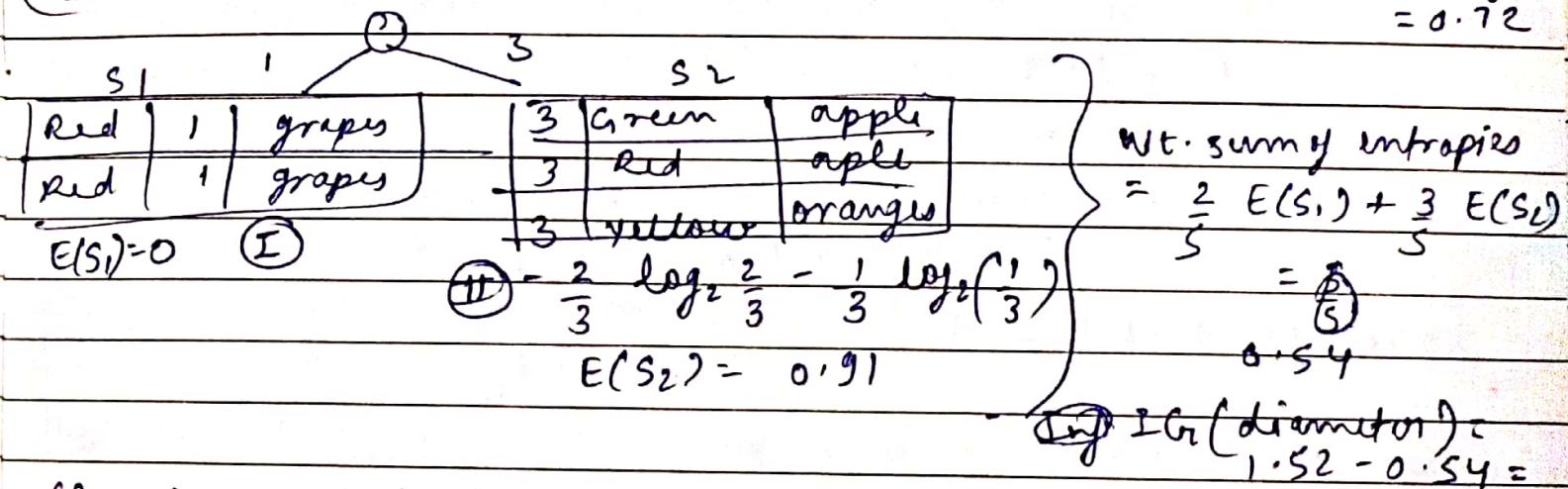
$$= 1.52$$



weighted sum
of entropies =

$$\frac{2}{5} \times E(S_1) + \frac{2}{5} \times E(S_2) + \frac{1}{5} \times E(S_3)$$
 $= 0.8$

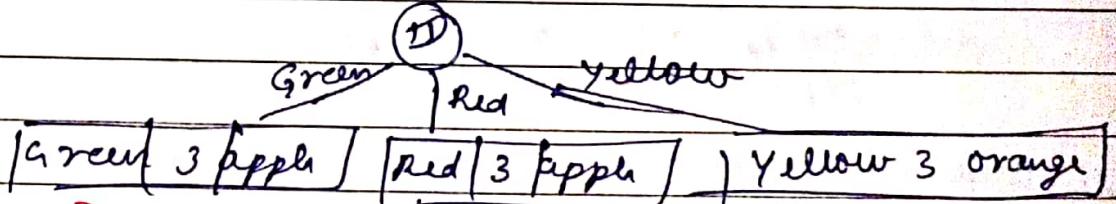
Information (color) = 1.52
Gains
 $= 0.8$
 $= 0.72$



$\text{IG}_r(\text{diameter}) = 1.52 - 0.54 = 0.98$

Clearly $\text{IG}_r(\text{diameter}) > \text{IG}_r(\text{color})$
Division on diameter

(I): Leaf node
Single Target
Class



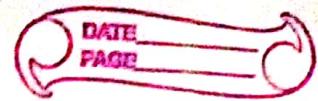
Problem with ID 3

① If scholar number is a feature then we will get all leaf nodes in single level. Highly biased for feature.

Attribute is having more number of ~~in~~ different instances.
It will dominate information gain (max. IG for such feature)

... of each
leaf node
corresponds to ...

28/02/24
Wednesday.



C4.5 Algorithm

- (1) It is an improvement over ID3
- (2) ID3 is biased towards attributes with more diff. values.
ex. Scholar number of student example.
- (3) To overcome this bias issue C4.5 to select attribute uses a purity measure "Gain Ratio" to identify the attribute with at least 2 or 3 values.
- (4) In C4.5 algorithm, the Information Gain measure used in ID3 algorithm for is normalized by considering another factor. "T-node" is a dataset at node.

$$\text{Gain Ratio}(A) = \frac{\text{Information Gain}(A)}{\text{split info}(T, A)}$$

$$\text{split info}(T, A) = - \sum_{i=1}^v \frac{|A_i|}{|T|} \log_2 \frac{|A_i|}{|T|}$$

where, the attribute A has got v distinct values $(a_1, a_2, a_3, \dots, a_v)$
containing a_i instances of a_i . In Entropy = instances of values.

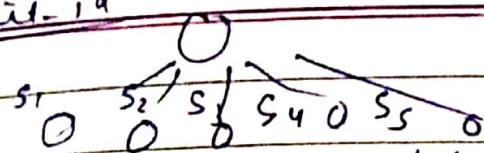
$$= -\frac{|a_1|}{|T|} \log_2 \frac{|a_1|}{|T|} - \frac{|a_2|}{|T|} \log_2 \frac{|a_2|}{|T|} \dots$$

Fruit-ID	Color	Diameter	Fruit
101	Green	3	apple
102	Red	3	apple
103	Yellow	1	grapes
104	Yellow	3	orange
105	Green	1	grapes

So: Entropy $E(S) = 1.52$

Calculate entropy, information gain, split info, gain ratio for all attributes.

~~fruit-id~~



$$E(S_1) = 0 \sim E(S_2) = \dots E(S_5) = 0$$

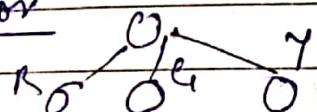
$$IG = 1.52 \\ (\text{fruit-id})$$

$$\text{Information Gain}(\text{fruit-id}) = 1.52 - \left(\frac{1}{5} \times 0 + \frac{1}{5} \times 0 \right) = 0$$

$$\text{Split info (fruit-id)} = - \left\{ \left(\frac{1}{5} * \log_2 \frac{1}{5} \right) \times 5 \right\} \\ \left(\frac{1}{5} \log_2 5 + \frac{1}{5} \log_2 5 \right) \approx 2.3219$$

$$\text{Gain ratio (fruit-id)} = \frac{1.52}{2.3219} = 0.656$$

color



$$\text{split ratio } \delta = \frac{2}{5} \log_2 \frac{2}{5} + \frac{2}{5} \log_2 \frac{2}{5} + \frac{1}{5} \log_2 \frac{1}{5}$$

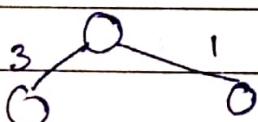
$$IG(\text{color}) = 0.72$$

$$E(\text{color}) = 0.8$$

$$\text{Gain ratio} = \frac{0.72}{0.8} = 0.9$$

→ best class

Diameter



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

$$= 0.97$$

$$IG(\text{diameter}) = 0.98$$

$$\text{Gain ratio} = \frac{0.98}{0.97} = 1.01$$

Clearly, ~~IG(diameter)~~

Gain ratio (diameter) is highest.

			diameter	
			2	3
101	Green	apple		
102	Red	apple		
104	Yellow	orange		
103	Red	grape		
105	Green	grape		

We were able to bypass fruit-id from selection ~~as~~ which was the problem in ID3.

64.02.2024
Monday

Gini-Index $\in [0, 0.5]$
clear cut

DATE
PAGE

most confusion

CART Algorithm

Classification and Regression Tree

↳ standard deviation instead of entropy

① It constructs the tree as a binary tree by recursively splitting a node into 2 nodes.

(If we have 3 classes then combine 2 classes so that we have 2 children only)

② Gini Index

$$\text{Gini Index} = 1 - \sum_{i=1}^k p_i^2$$

min value = 0

max value = 0.5

p_i = Probability that a data instance belongs to class C_i

$$\text{Gini Index} = 1 - \sum_{i=1}^k (P(C_i))^2$$

Entropy, Same

③ Every attribute is considered as a binary attribute which splits the nodes into 2 subsets S_1 & S_2

$$\text{Gini Index}(T, A) = \frac{|S_1|}{|T|} \text{Gini}(S_1) + \frac{|S_2|}{|T|} \text{Gini}(S_2)$$

Weighted sum of Gini Index

④ AGINI is calculated for each attribute

$$\Delta \text{Gini}(A) = \text{Gini}$$

$$\Delta \text{Gini}(A) = \text{Gini}(T) - \text{Gini}(T, A)$$

Attribute 'A' used to split the tree 'T'

If an attribute 'A' has 3 distinct values (a_1, a_2, a_3)

Then possible subsets are $\{a_1\}, \{a_1, a_2\}, \{a_1, a_3\}, \{a_2\}, \{a_2, a_3\}, \{a_3\}, \{a_3, a_1\}, \{a_3, a_2\}$

X Y

combination of division

$\{a_1\}$	$\{a_2, a_3\}$	$\rightarrow \text{GiniIndex}(T, A) \rightarrow C_1$
$\{a_2\}$	$\{a_1, a_3\}$	$\rightarrow \text{GiniIndex}(T, A) \rightarrow C_2$
$\{a_3\}$	$\{a_1, a_2\}$	$\rightarrow \text{GiniIndex}(T, A) \rightarrow C_3$
X	Y	3 subsets are possible

$\{a_1, a_2\} \{a_3, a_4\}$

For 4
 $\sum_{a_1, a_2} \{a_3, a_4\}$ ✓
 DATE
 PAGE

Select $\min(S_{G_1}, G_2, G_3)$ → to split the node.

In CART algorithm, we need to compute the best splitting attribute & the best split subset is in the chosen attribute.

Q) Color	Diameter	Fruit	Create Decision Tree using CART algorithm.
Green	3	Apple	
Red	3	Apple	
Red	1	Grapes	
Yellow	3	Orange	
Green	1	Grapes	

$$(P(\text{Apple}))^2 + P(\text{Grapes})^2 + P(\text{Orange})^2$$

$$\text{Sol: GiniIndex}(T) = 1 - \sum_{i=1}^3 p_i^2 = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{1}{5}\right)^2 - \left(\frac{2}{5}\right)^2 \\ = 1 - \frac{4}{25} - \frac{1}{25} - \frac{4}{25} = \frac{16}{25} = 0.64$$

Consider Color:-

color has 3 possible values,

So the combination of subsets for split.

(i) $\{G_1\}, \{R, Y\}$	0.6 0.6	(ii) $\{R\}, \{G_1, Y\}$	0.6 0.6
(ii) $\{R\}, \{G_1, Y\}$	0.6 0.6	R 3 Apple	G 3 Apple
(iii) $\{Y\}, \{G_1, R\}$.	0.6 0.4 ✓	R $\left(\frac{1}{5}, 1\right)$ Grapes	Y 3 Orange
(iv) $\{G_1\}, \{R, Y\}$	0.6 0.6	G 1 Grapes	

$$\text{Gini}(S_1) = \frac{1}{2} = 0.5$$

$$\text{Gini}(S_2) = \frac{2}{3} = 0.67$$

$$\text{Gini}(T_2, (ii)) = \frac{2}{5} \times \frac{1}{2} + \frac{3}{5} \times \frac{2}{3} = 0.6$$

Green 3 Apple Red 3 Apple
 Green 1 Grapes Red 1 Grapes

$$\text{Gini}(S_1) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = \frac{2}{4} = 0.5$$

$$\text{Gini}(S_2) = 1 - \left(\frac{2}{3}\right)^2 = \frac{1}{3}$$

$$\text{Gini}(T_2, (ii)) = \frac{2}{5} \times \frac{1}{2} + \frac{3}{5} \times \frac{2}{3} = \frac{2}{5} + \frac{2}{5} = \frac{4}{5} = 0.8$$

Weighted Gini

(iii) S_2 $\rightarrow S_{2,1}$

6	3	orange
G	1	(S ₁)

$$\text{Gini}(S_1) = 1 - 1^2 = 0$$

G	3	Apple
R	3	Apple
R	1	grapes
G	1	grapes
		(S ₂)

$$\text{Gini}(S_2) = 1 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Gini}(T, \text{(iii)}) = \frac{1}{5} \times 0 + \frac{2}{5} \times \frac{1}{2} = \frac{2}{5} = 0.4$$

$$\text{Gini}_{\text{color}}(T, \Delta \text{Gini}(\text{color})) = \text{Gini}(T) - \text{Gini}(T, \text{color}) \\ = 0.64 - 0.4$$

$$\Rightarrow \Delta \text{Gini}(\text{color}) = 0.24$$

Diameter
 $\{1\} \quad \{2\} \quad \{3\}$

G	Apple	green	3	Apple	
R	1	grapes	Red	3	Apple
G	1	grapes	Yellow	3	orange
	(S ₁)		(S ₂)		

$$\text{Gini}(S_1) = 1 - 1^2 = 0$$

$$\text{Gini}(S_2) = 1 - \left(\frac{1}{3}\right)^2 = \frac{2}{3} \\ = 1 - \frac{5}{9} = \frac{4}{9} \\ = 0.444$$

$$\text{weighted Gini} \rightarrow \text{Gini}(T, \text{Diameter}) = \frac{2}{5} \times 0 + \frac{3}{5} \times 0.444 \\ = 0.27$$

$$\Delta \text{Gini}(\text{Diameter}) = \text{Gini}(T) - \text{Gini}(T, \text{Diameter}) \\ = 0.64 - 0.27 = 0.37$$

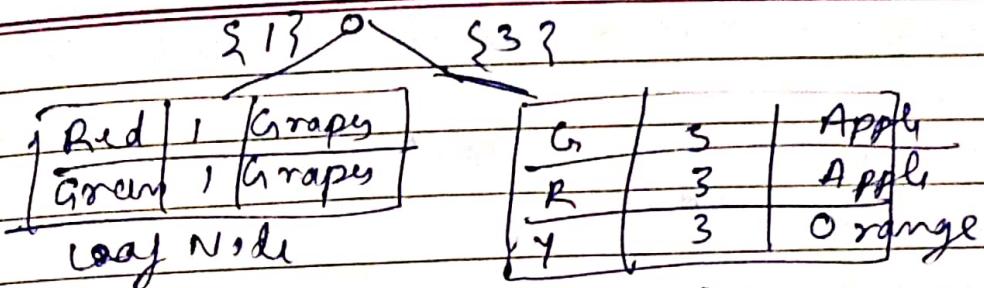
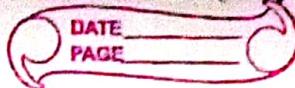
$$\Rightarrow \Delta \text{Gini}(\text{Diameter}) = 0.37$$

Clearly,

$$\Delta \text{Gini}(\text{Diameter}) > \Delta \text{Gini}(\text{color})$$

\Rightarrow Diameter is used to divide the node.

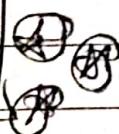
~~Some books, CART Algo. ⇒ No need for binary tree~~
 Some books, CART Algo. ⇒ No need for binary tree
 maybe that works better in their cost.



Same process ahead

Σ Y₃, Σ G, R₃

Naive Bayes give equal weightage to every attribute but Decision tree does not give equal attributes.



SUPPORT VECTOR MACHINE (SVM)

supervised learning algorithm.

classification

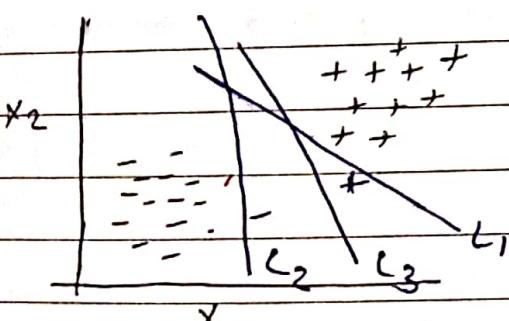
regression

Robust to outliers

work with non linear data as well.

Handling dimension of data.

Linear → line can separate the data.. into 2. halves.



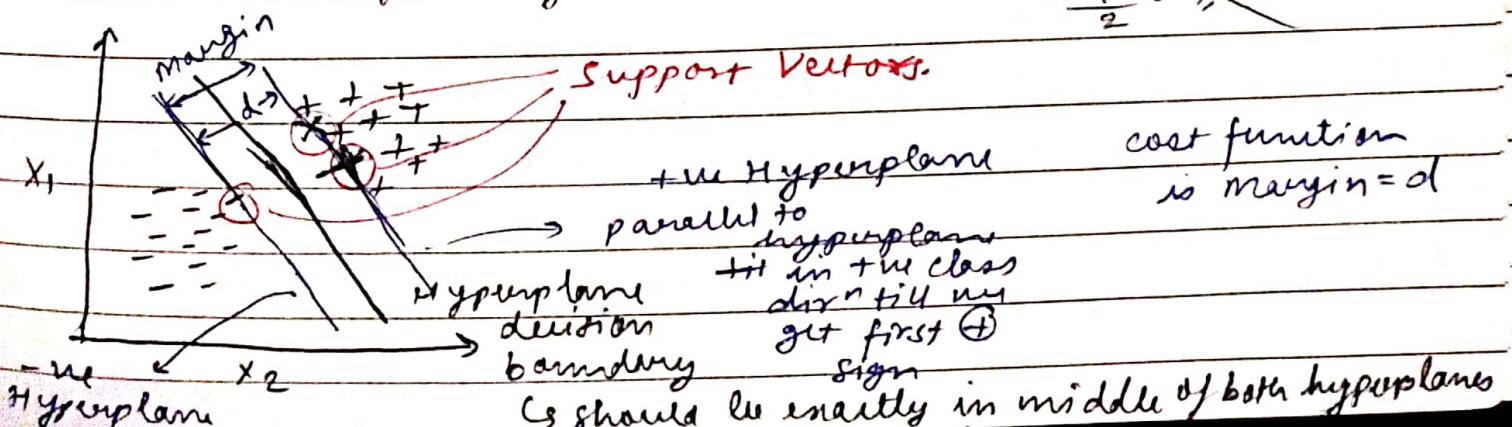
L₃ unable to classify ⊕ blue as ⊕ instead classify it as ⊖

Similar problem with L₂ ⊖ blue as ⊕

[k nearest neighbour would have given current answer]

Maximum margin hyperplane

$$\frac{d_1 + d_2}{2} = x$$



distⁿ b/w line & vector

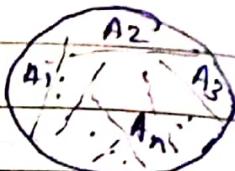
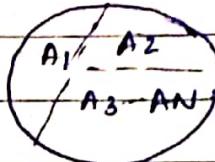
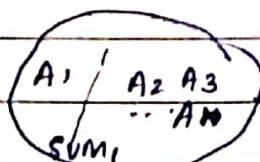
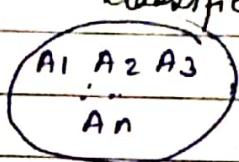
distⁿ two || line

magnitude of a vector

difference of 2 vectors

DATE
PAGE

SVM here is binary classification but can be extended for multiclass classification



A SVM for n Class classification.

at n class classification

The basic idea of SVM is distance b/w hyperplane & nearest datapoint should be maximum & hyperplane should avoid misclassification of data.

"Margin" is defined as the amount of space between the two classes as defined by the hyperplane.

The focus of SVM is to obtain "maximal margin classifier".

Types of SVM :-

- ① Hard margin SVM → Linearly separable data (previous page ex.)
- ② Soft margin SVM → Non-linearly separable data.

Support vectors :- Data points that fall on the supporting hyperplanes

① Hard Margin SVM as Binary classification :-

A SVM implements a binary classifier this means there are only 2 classes say +1 and -1.

This idea can be extended to multiclass SVM.

Let's consider a dataset with n-dimensions

i/p: $x_1, x_2, x_3, \dots, x_n$

o/p: $y = \{-1, +1\}$

Hyperplane is a n-dimensional generalization of a line, so

equation of Hyperplane is $w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + w_0 = 0$

$$5 + 3x_1 + 4x_2 = 0$$

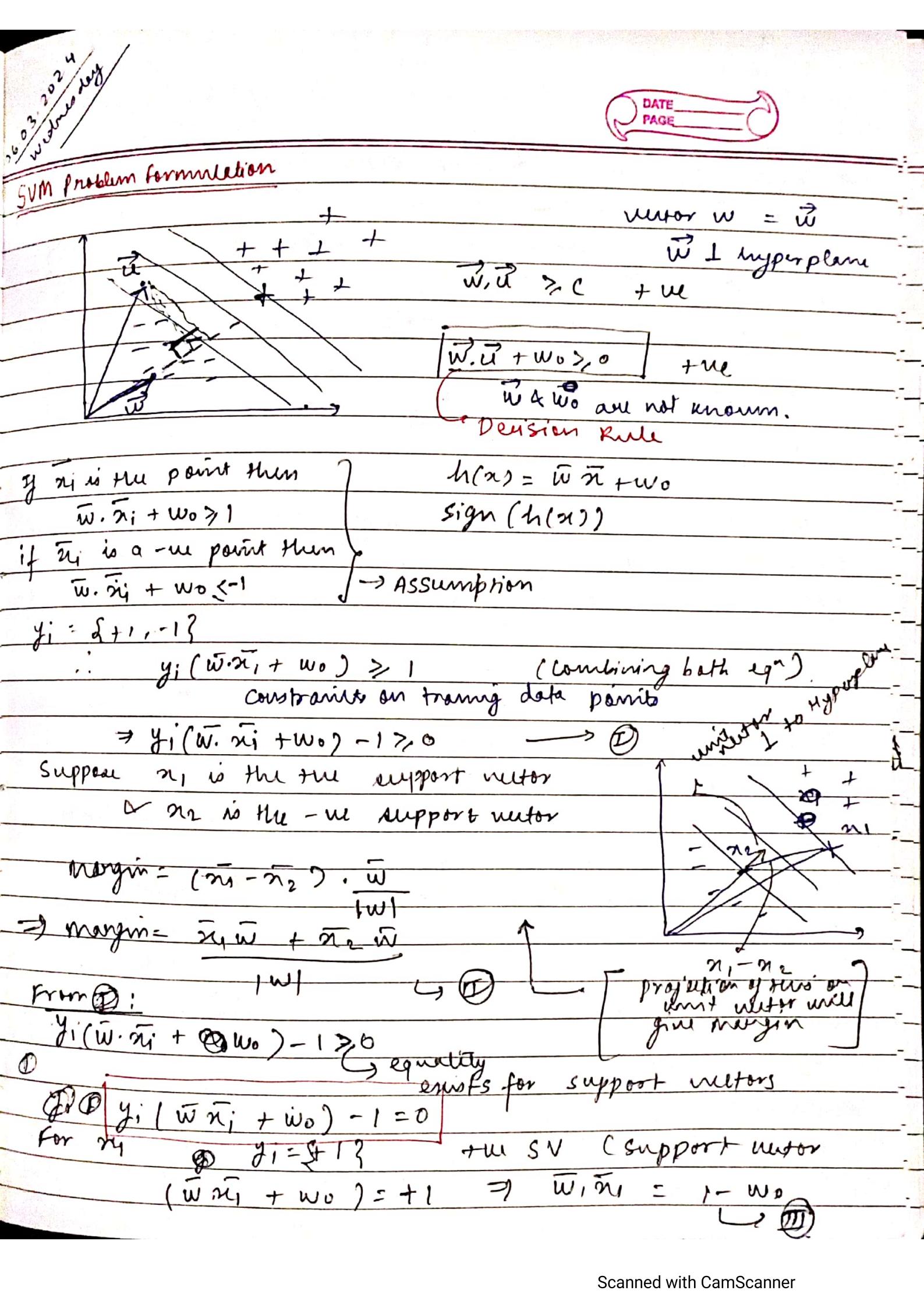
$$3 + 3x_1 + 4x_2 = 0$$

parallel lines

$$\Rightarrow \bar{w} \cdot \bar{x} + w_0 = 0$$

⊗ ⊗ ⊗

Assume $\left\{ \begin{array}{l} w_0 + w_1x_1 + w_2x_2 = 1 \quad (\text{top Hyperplane}) \\ w_0 + w_1x_1 + w_2x_2 = 0 \quad (\text{Hyperplane}) \\ w_0 + w_1x_1 + w_2x_2 = -1 \quad (-\text{bottom Hyperplane}) \end{array} \right.$



11.03.2024
Monday



$$\text{For } x_2, y_2 = -1 \quad \text{①} - w \cdot SV$$

$$-1 (\bar{w} \cdot \bar{x}_2 + w_0) = 1 \circ = 0$$

$$\Rightarrow -\bar{w} \cdot \bar{x}_2 = -1 + w_0 \rightarrow \text{④}$$

put ③ & ④ in ②:

$$\text{margin} = \frac{(1 - w_0) + (-1 + w_0)}{\|w\|}$$

$$\Rightarrow \boxed{\text{margin} = \frac{2}{\|w\|}} \quad \text{A} \quad \text{B}$$

max margin $\Rightarrow (1/w)$ minimize \Rightarrow minimize $\frac{1}{2} (1/w)^2$

Dual function \Rightarrow maximize \Rightarrow then original function minimizes - Optimization Technique

Hard margin SVM as an Optimization problem.

$$\underset{w, w_0}{\text{minimize}} \frac{1}{2} \|w\|^2$$

subject to the constraint

$$y_i (\bar{w} \cdot \bar{x}_i + w_0) - 1 \geq 0 \quad \forall x_i \in D$$

It's is an optimization problem with constraint.

Lagrangian Multipliers are one way of solving the optimization problem with respect to equality constraints.

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i (\bar{w} \cdot \bar{x}_i + w_0) - 1]$$

This is Lagrangian multipliers.

Actual f' is minimization

Maximizations \therefore we will minimize

α is Lagrangian Multiplier

$$\underset{\alpha}{\text{maximize}} \quad L(w, w_0, \alpha)$$

$$L = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i y_i \bar{w} \cdot \bar{x}_i - \sum_{i=1}^m \alpha_i y_i w_0 + \sum_{i=1}^m \alpha_i \rightarrow \text{I}$$

$$\frac{dL}{d\alpha_0} = 0 + 0 + \sum_{i=1}^m \alpha_i y_i + 0 = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0 \rightarrow \textcircled{I}$$

$$\frac{dL}{d\bar{w}} = \bar{w} - \sum_{i=1}^m \alpha_i y_i \bar{x}_i - 0 + 0 = 0 \\ \Rightarrow \bar{w} = \sum_{i=1}^m \alpha_i y_i \bar{x}_i \rightarrow \textcircled{II}$$

put the value of \textcircled{II} in \textcircled{I} :

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i y_i \bar{w} \cdot \bar{x}_i + w_0 \sum_{i=1}^m \alpha_i y_i + \sum_{i=1}^m \alpha_i$$

put \textcircled{I} :

$$\Rightarrow L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i y_i \bar{w} \cdot \bar{x}_i + \sum_{i=1}^m \alpha_i$$

Substitute \textcircled{II} :

$$\Rightarrow L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i y_i \left(\sum_{j=1}^m \alpha_j y_j \bar{x}_j \right) \bar{x}_i + \sum_{i=1}^m \alpha_i$$

$$= \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j + \sum_{i=1}^m \alpha_i$$

from \textcircled{I} :

$$= \frac{1}{2} \sum_{i=1}^m \alpha_i y_i \bar{x}_i \sum_{j=1}^m \alpha_j y_j \bar{x}_j - \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j + \sum_{i=1}^m \alpha_i$$

$$= \sum_{i=1}^m \alpha_i \sum_{j=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j - \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j + \sum_{i=1}^m \alpha_i$$

$$\Rightarrow L(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{x}_i \bar{x}_j$$

\oplus \ominus
 \star

$$\frac{dL}{d\alpha_1} = 0, \frac{dL}{d\alpha_2} = 0, \frac{dL}{d\alpha_3} = 0, \dots, \frac{dL}{d\alpha_m} = 0$$

m equations, m unknown variables [$=$ No. of data points]
can be solved for $\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_m$ (numerically)

for support vectors
 $d_i > 0$
 $d_i = 0$ for not support vectors.

DATE _____
 PAGE _____

$$\max(L(\alpha)) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \bar{n}_i \cdot \bar{n}_j$$

subject to $\alpha_i \geq 0 \quad i=1, 2, \dots, m$

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad (\text{from (ii)})$$

From (iii): $\bar{w} = \sum_{i=1}^m \alpha_i y_i \bar{n}_i$ ~~(prev)~~ ~~(curr)~~ ~~(curr)~~

Equality for support vectors:- (found earlier) (prev (curr))

$$y_i (\bar{w} \cdot \bar{n}_i + w_0) = 1$$

$$\Rightarrow y_i y_i (\bar{w} \cdot \bar{n}_i + w_0) = y_i \quad [\text{multiplied by } y_i]$$

$$\Rightarrow \bar{w} \cdot \bar{n}_i + w_0 = y_i$$

$$\Rightarrow w_0 = y_i - \bar{w} \cdot \bar{n}_i$$

$y_i = \pm 1, -1 \}$
 $y_i \cdot y_i = y_i^2$
 $\therefore y_i^2 = \pm 1 \}$
 $\therefore y_i^2 = +1 \}$

$$\Rightarrow w_0 = \frac{1}{S} \sum_{i=1}^S (y_i - \bar{w} \cdot \bar{n}_i)$$

Q) Three data point as support vectors

	x_1	x_2	y
1	1	2	-1
2	-1	2	-1
3	-1	-2	1

~~S.H.~~

$$L(\alpha) = \sum_{i=1}^3 \alpha_i - \frac{1}{2} \sum_{i=1}^3 \sum_{j=1}^3 \alpha_i \alpha_j y_i y_j \bar{n}_i \cdot \bar{n}_j$$

$$= \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} \sum_{i=1}^3 (\alpha_1 \alpha_1 y_1 y_1 \bar{n}_1 \cdot \bar{n}_1 + \alpha_1 \alpha_2 y_1 y_2 \bar{n}_1 \cdot \bar{n}_2 +$$

~~$\alpha_1 \alpha_3 y_1 y_3 \bar{n}_1 \cdot \bar{n}_3 + \alpha_2 \alpha_2 y_2 y_2 \bar{n}_2 \cdot \bar{n}_2 + \alpha_2 \alpha_3 y_2 y_3 \bar{n}_2 \cdot \bar{n}_3 + \alpha_3 \alpha_3 y_3 y_3 \bar{n}_3 \cdot \bar{n}_3$~~

$$= \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} [(\alpha_1 \alpha_1 y_1 y_1 \bar{n}_1 \cdot \bar{n}_1 + \alpha_2 \alpha_2 y_2 y_2 \bar{n}_2 \cdot \bar{n}_2 + \alpha_3 \alpha_3 y_3 y_3 \bar{n}_3 \cdot \bar{n}_3)$$

$$+ (\alpha_1 \alpha_2 y_1 y_2 \bar{n}_1 \cdot \bar{n}_2 + \alpha_2 \alpha_3 y_2 y_3 \bar{n}_2 \cdot \bar{n}_3 + \alpha_3 \alpha_1 y_3 y_1 \bar{n}_3 \cdot \bar{n}_1)]$$

$$+ \alpha_1 \alpha_3 y_1 y_3 \bar{n}_1 \cdot \bar{n}_3 + \alpha_2 \alpha_3 y_2 y_3 \bar{n}_2 \cdot \bar{n}_3 + \alpha_3 \alpha_1 y_3 y_1 \bar{n}_3 \cdot \bar{n}_1]$$

↪ 9 terms

29.2.2024
M.2

Transpose

$$y_i \cdot \bar{y}_i = (y_i \cdot \bar{y}_i) + (y_i \cdot \bar{y}_i)^T$$

DATE _____
PAGE _____

$$n = \sum_{i=1}^m y_i \cdot \bar{y}_i \cdot \bar{x}_i \cdot x_i$$

Diagram showing the calculation of the inner product $y_i \cdot \bar{y}_i$ and the outer product $(y_i \cdot \bar{y}_i)^T$ for three vectors x_1, x_2, x_3 . The vectors are represented as columns of numbers. The inner product is calculated by summing the products of corresponding elements from y_i and \bar{y}_i . The outer product is calculated by multiplying $y_i \cdot \bar{y}_i$ with each column of x_i .

$$\begin{array}{c|ccc} y & x_1 & x_2 & x_3 \\ \hline -1 & 1 & 2 \\ -1 & -1 & 2 \\ -1 & -1 & 2 \\ 1 & -1 & -2 \end{array}$$

$$\begin{matrix} j \rightarrow 1 & 2 & 3 \\ i \downarrow & x_1 x_2 & x_2 x_3 \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \left[\begin{matrix} -5 & 3 & 5 \\ 3 & 5 & 3 \\ 5 & 3 & 5 \end{matrix} \right] \end{matrix}$$

→ main

Representation of complex calculation

$$\sum_{i=1}^m \sum_{j=1}^m y_i \cdot \bar{y}_i \cdot x_i \cdot x_j$$

$$L(\alpha) = \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} [5\alpha_1 \alpha_2 + 3\alpha_1 \alpha_2 + 5\alpha_1 \alpha_3 + 3\alpha_2 \alpha_1 + 5\alpha_2 \alpha_2 + 3\alpha_2 \alpha_3 + 5\alpha_3 \alpha_1 + 3\alpha_3 \alpha_2 + 5\alpha_3 \alpha_3]$$

$$L(\alpha) = \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} [5\alpha_1^2 + 3\alpha_1 \alpha_2 + 5\alpha_1 \alpha_3 + 3\alpha_2 \alpha_1 + 5\alpha_2^2 + 3\alpha_2 \alpha_3 + 5\alpha_3 \alpha_1 + 3\alpha_3^2 + 3\alpha_3 \alpha_2]$$

$$\begin{aligned} \frac{dL}{d\alpha_1} &= 1 - \frac{1}{2} (10\alpha_1 + 3\alpha_2 + 5\alpha_3 + 3\alpha_2 + 5\alpha_3) = 0 \\ &= 1 - \frac{1}{2} (10\alpha_1 + 6\alpha_2 + 10\alpha_3) = 0 \\ &\Rightarrow 5\alpha_1 + 3\alpha_2 + 5\alpha_3 = 1 \quad \rightarrow \textcircled{I} \end{aligned}$$

$$\begin{aligned} \frac{dL}{d\alpha_2} &= 1 - \frac{1}{2} (3\alpha_1 + 3\alpha_1 + 10\alpha_2 + 3\alpha_3 + 3\alpha_3) = 0 \\ &= 1 - \frac{1}{2} (3\alpha_1 + 5\alpha_2 + 3\alpha_3) = 0 \quad \rightarrow \textcircled{II} \end{aligned}$$

$$\begin{aligned} \frac{dL}{d\alpha_3} &= 1 - \frac{1}{2} (-5\alpha_1 + 3\alpha_2 + 5\alpha_1 + 3\alpha_2 + 10\alpha_3) = 0 \\ &\Rightarrow 5\alpha_1 + 3\alpha_2 + 10\alpha_3 = 1 \quad \rightarrow \textcircled{III} \end{aligned}$$

$$\textcircled{I} - \textcircled{III}: 15\alpha_1 + 9\alpha_2 + 15\alpha_3 = 3$$

$$15\alpha_1 + 25\alpha_2 + 15\alpha_3 = 5$$

$$16\alpha_2 = 2 \quad \Rightarrow \alpha_2 = \frac{1}{8}$$

$$5\alpha_1 + 5\alpha_3 = 1 - \frac{3}{8} = \frac{5}{8}$$

$$\Rightarrow \alpha_1 + \alpha_3 = \frac{1}{8}$$

○ kept optimization
 ○ dual satisfactory
 cond¹
 DATE _____
 PAGE _____

$$\alpha_1 + \alpha_3 = \frac{1}{8}$$

$$\left[\begin{array}{l} \alpha_2 = \frac{1}{8} \\ \alpha_1 = 0 \end{array} \right]$$

point hyperplane
y = 0

α_2 & α_1 are related to -ve hyperplane

∴ we can make $\left[\begin{array}{l} \alpha_1 = 0 \\ \alpha_3 = \frac{1}{8} \end{array} \right]$

$$\therefore \left[\begin{array}{l} \alpha_3 = \frac{1}{8} \end{array} \right] \quad (+ve \text{ hyperplane})$$

$$\bar{w} = \sum_{i=1}^3 \alpha_i y_i x_i = 0 + \frac{1}{8} (-1)$$

$$= \alpha_1 y_1 x_1 + \alpha_2 y_2 x_2 + \alpha_3 y_3 x_3$$

$$= 0 + \frac{1}{8} (-1) [-1 \ 2] + \frac{1}{8} (1) [-1 \ 2]$$

↳ vector

$$= \frac{1}{8} [1 \ -2] + \frac{1}{8} [-1 \ -2]$$

$$\Rightarrow \left[\begin{array}{l} \bar{w} = \left(\begin{array}{l} w_0 \\ w_1 \\ w_2 \end{array} \right) = \frac{1}{8} \left(\begin{array}{l} 0 \\ 0 \ 1 \\ -1 \ 2 \end{array} \right) \end{array} \right] = \left[\begin{array}{l} 0 \\ -\frac{1}{8} \\ -\frac{1}{4} \end{array} \right]$$

$w_1 = 0$
 $y_2 = 1$
 $y_3 = -1$

$$w_0 = \frac{1}{8} \sum_{i=1}^3 (y_i - \bar{w} \bar{x}_i)$$

evaluate it for support vector 1 & 2.

$$= \frac{1}{2} (y_2 - \bar{w} \bar{x}_2 + y_3 - \bar{w} \bar{x}_3)$$

$$= \frac{1}{2} (-1 - [0 \ -\frac{1}{2}] \cdot [-1 \ 2] + 1 - [0 \ -\frac{1}{2}] \cdot [-1 \ -2])$$

$$= \frac{1}{2} (-1 + 1 + 1) = 0 \Rightarrow \boxed{w_0 = 0}$$

$$\bar{w} = \left[\begin{array}{l} 0 \\ -\frac{1}{8} \\ -\frac{1}{4} \end{array} \right]$$

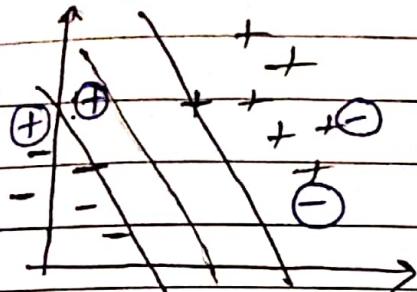
$$w_1 x_1 + w_2 x_2 + w_0 = 0$$

$$\Rightarrow 0x_1 + 0 - \frac{1}{2} x_2 + 0 = 0$$

$$\Rightarrow \boxed{x_2 = 0} \quad \text{Final eqn}$$

Hard margin SVM → rigid to outliers
 Flexibility is required → soft margin SVM

Soft Margin SVM



All +ve should be on one side & all -ve sides on other

not rigid on this case

- ① Hard margin SVM is very rigid for 2 constraints of SVM:
 - i) Decision boundary should be as far as from the nearest data point.
 - ii) Avoid misclassification.
- ② But the real world problems are complex and often datasets are affected by noise & outliers.
[data will not be completely linearly separable]
- ③ But dataset is almost linearly separable, (relaxation in margin, Balance b/w margin & misclassification)

$$\min_{\mathbf{w}_1, \mathbf{w}_0} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

distance error classification error



pronunciation:

subject to constraint $\xi_i \geq 0$ & y_i

$\xi_i \Rightarrow$ slack variable

y_i correctly classified
 $\xi_i = 0$

Misclassification
 $\xi_i > 0$

$$y_i (\mathbf{w} \cdot \mathbf{x}_i + w_0) \geq 1 - \xi_i$$

you can tune it

C ad more
 ξ_i more important to
 C is less
misclassification error

distance error \rightarrow Margin error importance

Lagrangian dual of this is similar to Hard margin SVM
only change $\alpha \geq 0$ & $\alpha < c$

E

O

S₁

a

A