

Supervised learning

- Discrete output (Classification)
- Real-valued (Regression).

$x_1, x_2, \dots, x_D \rightarrow$ features

$y_1, y_2, \dots, y_K \rightarrow$ outputs

$$x_i \in \mathbb{R}^D$$

Classification

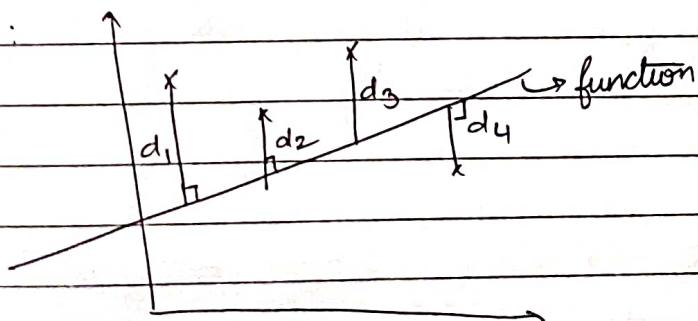
- Binary
- Multi-class.

Regression

Single-output regression : $f(x_i) \rightarrow y$.

linear — $y = mx + c$.

e.g. :



$$\min \sum_{i=1}^N d_i^p$$

$d_i \rightarrow$ euclidean distance.

Minkowski's Distance (L^p norm)

$$\sqrt[p]{|x_2 - x_1|^p + |y_2 - y_1|^p}$$

$p=1 \rightarrow L^1$ norm

$p=2 \rightarrow$ euclidean distance

to account
for error \rightarrow minimize $\sum (\hat{y}_i - y_i)^2$

residual
error

actual
predicted

Here,

loss function,

$$\min_{m, c} \sum_{i=1}^N (cmx_i + c - y_i)^2$$

\hat{y}_i y_i

$f(m, c) = \min \{ \text{sum of squared residual errors} \}$.

Multivariate Linear Regression,

$$\hat{y} = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_D x_D$$

α_0 intercept

$$\text{Roots} \rightarrow \frac{\partial J}{\partial \alpha_0} = \frac{\partial J}{\partial \alpha_1} = \dots = \frac{\partial J}{\partial \alpha_D} = 0.$$

Correlation

$$\frac{N \sum_{i=1}^N (a_i - \bar{a})(b_i - \bar{b})}{\text{std}(A) \text{std}(B)} = \frac{\sqrt{\sum (a_i - \bar{a})^2}}{\sqrt{\sum (b_i - \bar{b})^2}}$$

$$\sigma^2(A) = \frac{1}{N} \sum (a_i - \bar{a})^2$$

$$\sigma_A = \sqrt{\sigma^2(A)}$$

$$= \frac{\sum (a_i - \bar{a})(b_i - \bar{b})}{\sqrt{\sum (a_i - \bar{a})^2} \sqrt{\sum (b_i - \bar{b})^2}}$$

$$\text{Covariance} = \frac{1}{N} \sum (a_i - \bar{a})(b_i - \bar{b})$$

* Correlation $\in [-1, 1]$.

Covariance $\in (-\infty, \infty)$.

Linear Regression for Classification

$$\hat{y}_i = \alpha_0 + \alpha_1 x_1^i + \alpha_2 x_2^i + \dots + \alpha_D x_D^i$$

$i \rightarrow i^{\text{th}} \text{ row}$.

$$\hat{y}_i \in (-\infty, \infty).$$

Q. why do we take $\mathcal{E}_0, 13$ & not $\mathcal{E}_1, 23$ for binary classification?

It results in giving more weightage to higher class labels.
in order to minimize difference w/ separator.

How to use Linear Regression w/ Multi-class Dataset

1) One vs one.



(we get 2 comb.)

2) One vs all



(we get n comb.)

3) Error correcting output code.

One vs One

Take 1 1 & find max. of all sq' colⁿs.

One vs All

Find max ~~vectorizing row~~ w/most.

One v One has more data loss than One v All

* In 1 v 1, only $\frac{N}{K} \times Q$ samples are taken b/c $N - \frac{2N}{K}$ are lost

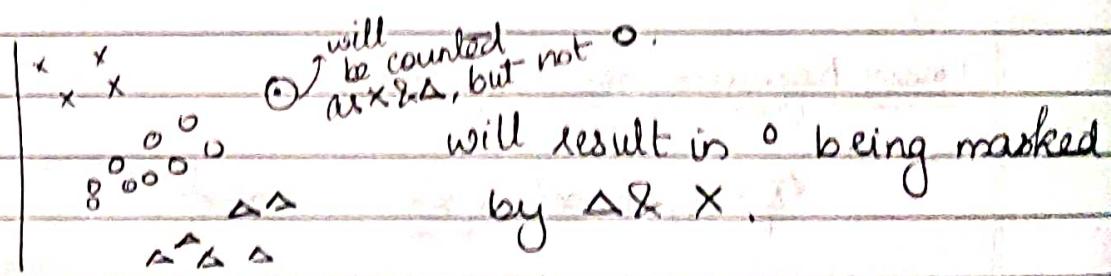
but in 1 v All, all are taken.

v. high if
K is high

But, if many categories,

1 v All has way more 0 than 1 a.k.a.
So, both are inaccurate if categories errors
are a lot.

Also,



Sigmoid $f^n \rightarrow \frac{1}{1+e^{-x}}$ always returns $\in [0, 1]$.

Error Correcting Output Codes

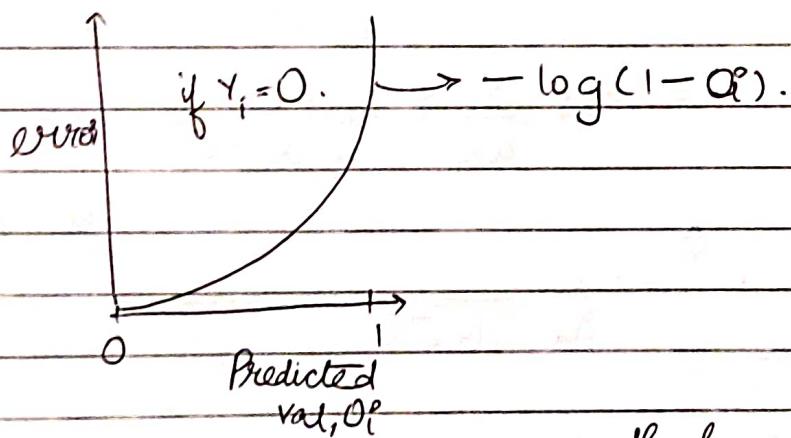
(Check later).

Logistic Regression (Binary Classifier). [here, it is better to use 1v1].

$$\hat{y}_i = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n \\ = \sum \alpha_j x_j \text{ if } x_0 = 1.$$

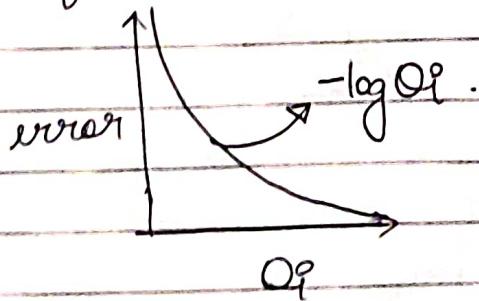
To convert vals b/w $(-\infty, \infty)$, we can use sigmoid to obtain $[0, 1]$ range.

Since sq. loss fn won't return sq. fn for e^{-x} , we have a non convex fn & need a convex optimisation method.



* If $f(x)$ is increasing, $\log(f(x))$ is increasing.

If $Y_i = 1$



Term 1 $\rightarrow Y_i = 1$

$$\sum [-y_i \log O_i] + [\textcircled{2} -(1-y_i) \log(1-O_i)]$$

$$\text{error} = \sum^k [-y_i \log O_i - (1-y_i) \log(1-O_i)]$$

①

Now,

$$\hat{O}_i = \frac{1}{1 + e^{-\hat{y}_i}}, \quad \hat{y}_i = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_D x_D.$$

Here,

$$\alpha_i^{t+1} = \alpha_i^t - \gamma \cdot \frac{\partial J}{\partial \alpha_i}, \quad \text{up}$$

Consider term ①:

$$\sum \frac{\partial}{\partial \alpha_i} (-y_i \log O_i) = -y_i \sum \frac{\partial \log O_i}{\partial \alpha_i}.$$

$$= -y_i \frac{\partial}{\partial \alpha_i} \sum \log \left(\frac{1}{1 + e^{-\hat{y}_i}} \right).$$

$$= -y_i \sum \underbrace{(1 + e^{-\hat{y}_i})}_{O_i} \times \frac{\partial}{\partial \alpha_i} \frac{1}{1 + e^{-\hat{y}_i}}$$

$$= -y_i \sum \frac{1}{O_i} \times \left(\frac{1}{e^{-\hat{y}_i} + 1} \right)^2 \times \frac{\partial}{\partial \alpha_i} (e^{-\hat{y}_i} + 1)$$

$$= -y_i \sum \left[\frac{1}{O_i} \times \frac{1}{(1 + e^{-\hat{y}_i})^2} \times e^{-\hat{y}_i} \times \frac{\partial}{\partial \alpha_i} \hat{y}_i \right].$$

From ①

$$= y_i \sum \left[O_i \times e^{-\hat{y}_i} \frac{\partial}{\partial \alpha_i} \hat{y}_i \right].$$

↓
 \hat{y}_i
 $\rightarrow x_i$

$$\hat{y}_i - 1$$

$$\text{Term ①: } \sum y_i \cdot \alpha_i^t \cdot \left(\frac{1}{\alpha_i^t} - 1 \right) x_i^t$$

$$\sum y_i \cdot \alpha_i^t \cdot \left(\frac{1}{\alpha_i^t} - 1 \right) x_i^t$$

$$\boxed{\sum y_i \cdot \alpha_i^t \left(\frac{1}{\alpha_i^t} - 1 \right) x_i^t = \sum y_i \cdot \alpha_i^t \left(\frac{1 - \alpha_i^t}{\alpha_i^t} \right) x_i^t}$$

$$\boxed{\text{① } \hat{y}_i^t = \sum y_i x_i (1 - \alpha_i^t)}$$

On simplifying & adding term ②,
we get:

~~$$\alpha_{i+1}^t = \alpha_i^t - \lambda \cdot \sum [\alpha_i^t - \hat{y}_i^t] x_i^t$$~~

Compare this w/ gradient descent for
linear regression,

$$\alpha_{i+1}^t = \alpha_i^t - \lambda \cdot \sum (\hat{y}_i^t - y_i) x_i^t$$

$\lambda \rightarrow$ step length.

* Here, decision boundary does NOT change.

α_i^t can also be represented by

$$P(Y=1 | x_1=3, \dots, x_d=5, \bar{y})$$

for eg.

Error Correcting Output Codes (ECOC).

Sample length is chosen b/w K & K_{C_2} .

can go

beyond this but it is
not that useful.

$$K \leq l \leq K_{C_2}$$

Class	B ₁	B ₂	B ₃	B ₄	B ₅	B ₆	
1	1	0	1	0	1	1	represents class 1.
2	1	1	1	1	1	1	
3	1	1	1	1	1	1	
4	1	1	1	1	1	1	
5	1	1	1	1	1	1	

random binary no.

random ECOC matrix

To verify that the binary labels are diff:

do XOR & count how many 1's are there.

The strings should be far apart by a length of $\frac{l}{2}$ minimum

Horizontal separability

Vertical Separability

No. of 1's in a col. should approx be equal.

To classify, generate outcome string & do XOR, and then put it in the class w/ least no of 1's.

Overfitting:

If 2 models M_1 & M_2 ,

When,

training error

$$\text{tr}(M_1) < \text{tr}(M_2)$$

But test error,

$$\text{test}(M_1) > \text{test}(M_2).$$

M_1 is overfitting.

Regularizations (to avoid overfitting):

loss f^n ,

$$\left[\sum -y_i \log \theta_i - (1-y_i) \log(1-\theta_i) \right] + \lambda \times \text{Penalty term}$$

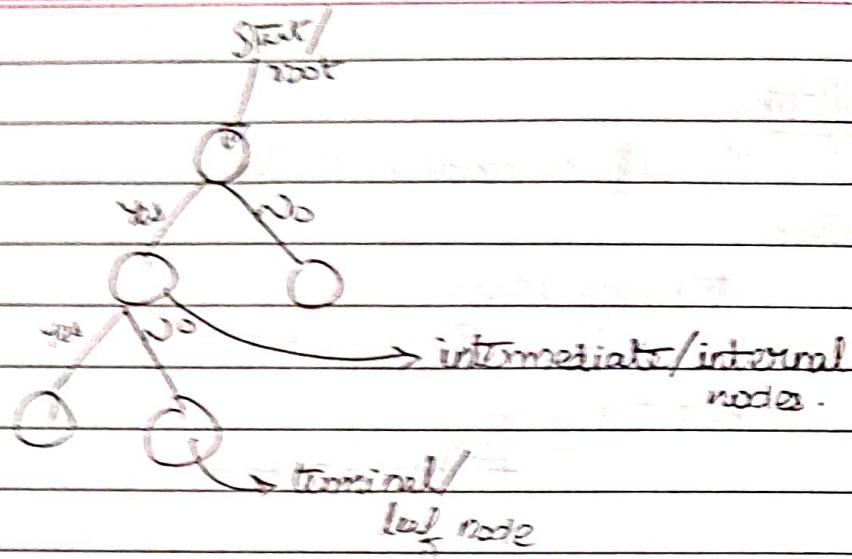
same regularisation
as linear.

$$L_1 \rightarrow \sum_{j=1}^D |\alpha_j|$$

$$L_2 \rightarrow \sum_{j=1}^D \alpha_j^2$$

$$\text{Elastic net} \rightarrow \lambda_1 \times L_1 \text{ Penalty} + \lambda_2 \times L_2 \text{ Penalty}$$

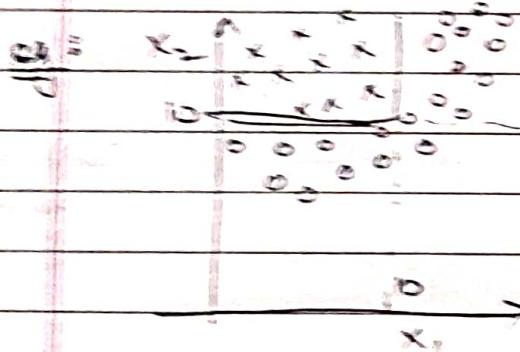
$$\lambda = \frac{\lambda_1}{\lambda_1 + \lambda_2}$$

Decision Tree

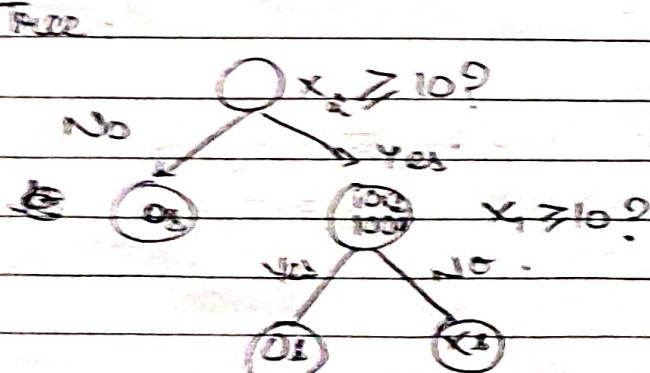
Root is a internal node conditions

Leaf nodes have no conditions, they make a decision.

- Can directly be used for multi-class datasets.

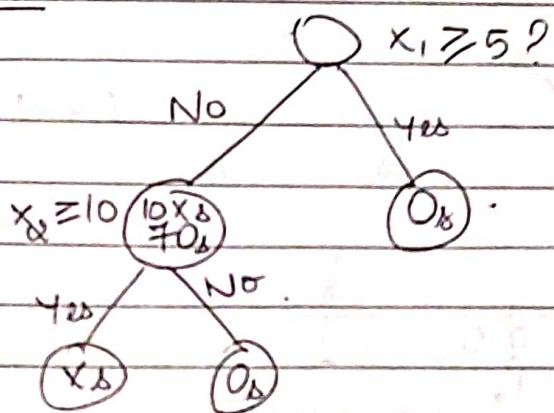


One separator line is impossible.
So, here there can be two lines to
help w/ this.



The type of decision boundary is called piecewise linear

alternative tree :



Clearly, 2nd one has a more pure root partition than 1st tree.
(10:7 vs 10:10). So we would prefer this.

How we test this → algorithm.

How to decide Purity

eg:	<table border="1"> <tr> <td>x</td><td>x₁</td></tr> <tr> <td>y</td><td>O₈</td></tr> </table>	x	x ₁	y	O ₈
x	x ₁				
y	O ₈				

Purest → One is 0

Most impure → x = y.

Shannon's entropy → $\sum_{i=1}^c -p_i \log p_i$, c → no. of classes.
p → probability

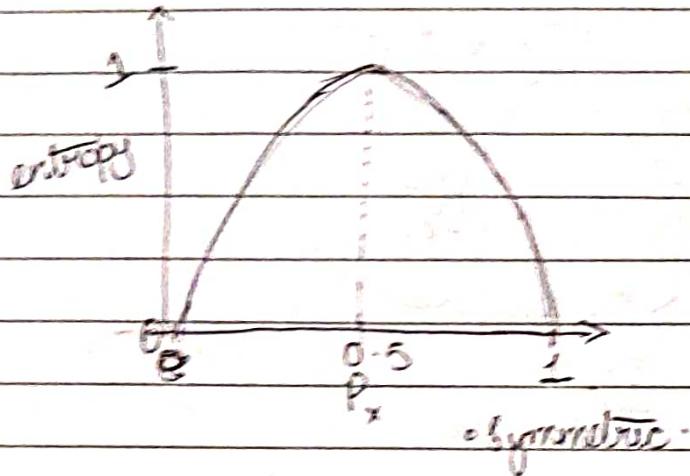
For eg,

$$5 \text{ } x_8 \& 5 O_8 \rightarrow -\frac{1}{2} \log \frac{1}{2} \times 2 = -\log_2 2^{-1} = +1$$

(Most impure case).

$$10x_8 \& O_8 \rightarrow -\frac{1}{2} \log \frac{1}{2} = 0$$

this entropy should be minimised.



(only for 2 classes)

for 3 classes, it is $> \frac{1}{2}$ in max. case.

Entropy would increase w/ no. of classes.

Gain in purity $\hat{=} \text{entropy of parent} - \sum_{\text{weighted}} \text{entropy for each child}$

Information gain $\hat{=} \text{entropy of parent} - w_1 \sum \text{entropy of each child}$.

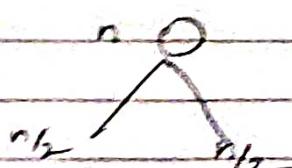
$$w_1 = \frac{\text{no. of samples the child uses}}{\text{no. of samples the parent uses}}$$

e.g.: 30 samples

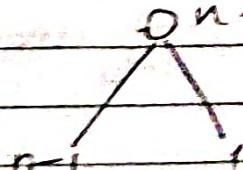
13 samples
 $w_1 = 13/30$

17 samples left
 $w_1 = 17/30$

No. of partitions for categorical features = no. of categories.
 To avoid overfitting & having partitions w/ $(n-1) \times 1$ sample,
 we take threshold as mean.



is better than



We can also take ranges of diff. ages (buckets) to create multiple children if each bucket has equal no. of samples.

- * Categorical features typically have higher entropy reduction (info. gain) than numerical features.

Other metrics (learn for exams)

- Gini impurity index
- Gain ratio.
- Cost sensitive ratio

Greedy Algorithm → only checks best choice at every step.
Algorithm

- * Recursively partitions data into pure, labelled partitions

Dataset (in slides)

$$\text{entropy}_{\text{root}} = \left(-\frac{9}{14} \log \frac{9}{14} - \frac{5}{14} \log \frac{5}{14} \right) \quad \text{for buys computer}$$

Consider all 4 features $= 0.9403$.

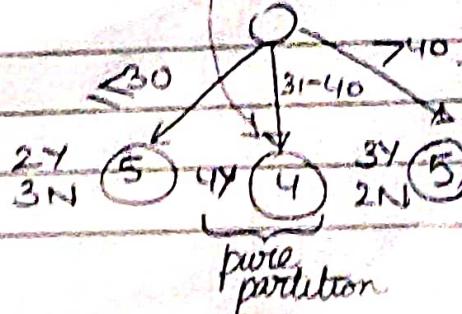
For age \rightarrow low & premium \rightarrow high

$$\text{info. gain} = \text{entropy}_{\text{root}} - \text{entropy}_{\text{left}} - \text{entropy}_{\text{right}}$$

For age \rightarrow entropy

$$\text{entropy}_{\text{left}} = \frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5} = \frac{2}{5} \log 2/5 + 0$$

$$\text{info. gain} = -2 \left[\frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5} \right].$$



$$\text{Info. gain} = \text{parent entropy} - \left[\frac{5}{14} \left(\frac{1}{2} E_{C_1} \right) + \frac{5}{14} E_{C_3} + 0 \right]$$

$$= 1.11 - \frac{5}{7} E_{C_1/C_2}$$

$$= 0.9403 - \frac{5}{7} \times 0.9710$$

$$= 0.2467$$

Similarly, for student,

$$IG = 0.1519$$

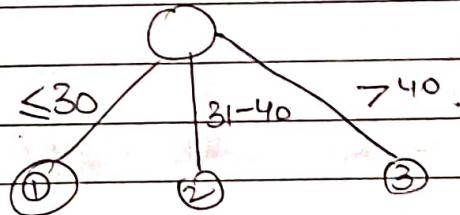
income,

$$IG = 0.0292$$

C.R.,

$$IG = 0.0481$$

We choose the highest IG & partition by age.

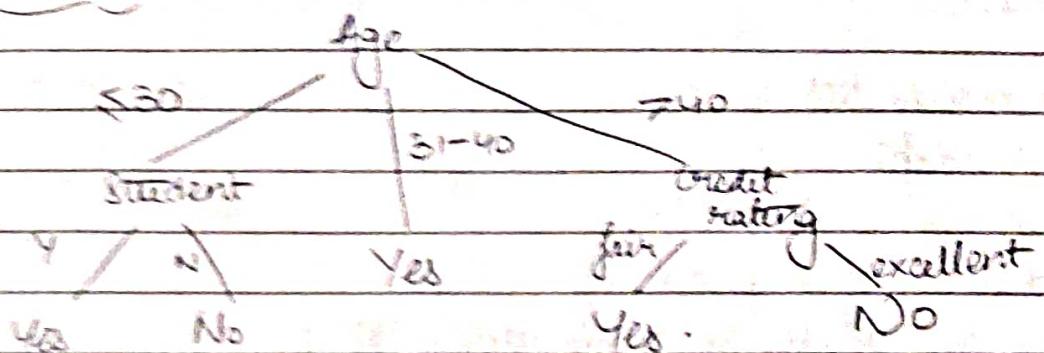


Look at node ①.

5 samples

3Y
2N

$$IG = 0.9710 - (\text{for each feature})$$

Final Tree

→ If we increase tree depth, it is likely to overfit. ^{disadvantage}

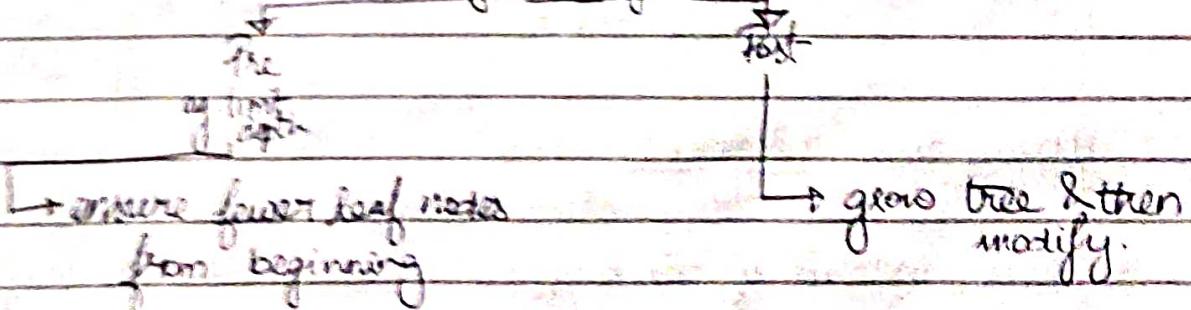
Randomised Trees

At every step take the top 'k' features & then choose best one ^{we choose}

→ It reduces chance of overfitting.

Other ways to avoid overfitting

1. Limit tree depth
2. Use pruning strategies or other approaches.

Breeding strategies

Pre-pruning

- 1) Randomised tree
- 2) Limit depth
- 3) Limit leaf nodes.
- 4) Put threshold on probability
 - 5) If v. few samples in node, don't split.
 - 6) If splitting results in v. few in 1 node & many in other, don't split

Post pruning

- 1) Merge redundant subtrees (if merging increases accuracy).

Subtree error pruning.

most effective strategy b/c we don't assume any parameters from the start.

↓
But, most computer intensive

- 2) Rule post-pruning (not v. effective).

we try to combine if-else conditions.

Decision trees for Regression

Since val. can't be same for regression here, instead of checking all in same class for leaf nodes, we check if they are close in value.

The metric we use to compute how diverse they are is standard deviation.

Support Vector Machines (SVM)

$$\Rightarrow \text{Minimise} \sum_{i=1}^N \sum_{j=1}^N (x_i - x_j)^2, \quad i < j, i \neq j.$$

$$= N \sum_{i=1}^N (x_i - \bar{x})^2. \quad \text{This lowers it from } O(n^2) \text{ to } O(n).$$

$$\Rightarrow \text{Minimise} \sum (x_i - \bar{x})^2$$

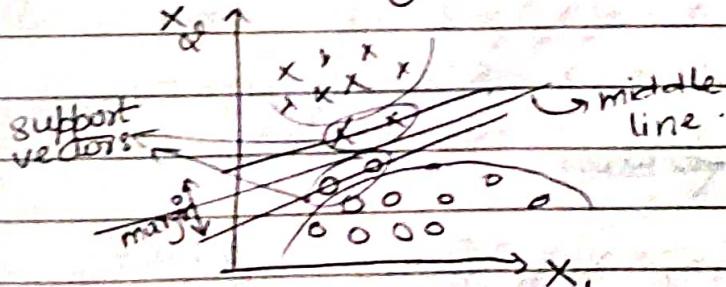
$$\Rightarrow \text{Minimise} \frac{1}{N} \sum (x_i - \bar{x})^2$$

$$\Rightarrow \text{Minimise} \sigma^2 \text{ (variance).}$$

+ We branch by picking the feature w_1 gives you the least variance on branching.
(we can also choose other metrics)

Support Vector Machines (SVM) (v. matrx)

Unlike logistic regression, we give more imp to pts closer to line.
To draw boundary,



we choose parallel lines' combo based on maximum margin.
SVM is also called max margin classifier.

The || line eqⁿ,

$$\alpha_1 x_1 + \alpha_2 x_2 = \text{constant}$$

both lines will have same eqⁿ w/ diff. constants.

For mathematical simplicity,

let const 1 = -const 2.

→ middle line = 0.

Eqs:

$$\left(\frac{\alpha_1}{K_0} \right) x_1 + \left(\frac{\alpha_2}{K_0} \right) x_2 = 1 ; \left(\frac{\alpha_1}{K_0} \right) x_1 + \left(\frac{\alpha_2}{K_0} \right) x_2 = -1 .$$

new α 's.

Middle line → $\alpha_1 x_1 + \alpha_2 x_2 = 0$

$$\rightarrow \alpha_1 x_1 + \alpha_2 x_2 = 1 ; \alpha_1 x_1 + \alpha_2 x_2 = -1 .$$

dist. b/w the lines, $d = |c_1 - c_2| / \sqrt{\alpha_1^2 + \alpha_2^2}$

$$\sqrt{\alpha_1^2 + \alpha_2^2} \quad \sqrt{\alpha_1^2 + \alpha_2^2}$$

⇒

Maximise d

$$\text{s.t. } (\alpha_1 x_1 + \alpha_2 x_2) \geq 1 \text{ if } y_i = 1$$

$$\alpha_1 x_1 + \alpha_2 x_2 \leq -1 \text{ if } y_i = -1 .$$

↓ becomes

$$y_i (\alpha_1 x_1 + \alpha_2 x_2) \geq 1 . \quad \text{3 hard SVM.}$$

We solve this using Lagrange multipliers.

Soft SVM,

$$\geq 1 - \epsilon$$

error

(allow some misclassifications to max. margin).

We maximise $\frac{1}{2} \alpha^T \alpha - \frac{1}{2} \alpha^T (\|w\|_2^2)$

is same as maximising | minimise

primal $\left\{ \begin{array}{l} \text{eqn} \\ \text{egn} \end{array} \right\} \frac{1}{2} \alpha^T \alpha - \frac{1}{2} \alpha^T (\|w\|_2^2)$

dual $\left\{ \begin{array}{l} \text{eqn} \\ \text{egn} \end{array} \right\} \frac{1}{2} \|w\|_2^2$

No. of constraints = N (no. of samples).

No. of parameters $\rightarrow N + d$

\hookrightarrow no. of features

Lagrange Multipliers

$\min f(x) \text{ s.t } g(x)$
becomes

$$f(x) + \lambda g(x)$$

where $\lambda > 0$.

* If $g(x)$ has \geq then $\lambda \geq 0$

$g(x)$ has $<$ then $\lambda < 0$.

* In regression, we take support vectors as OUTERMOST pts.



Ensemble Learning

Using multiple models based on sample.
then choose a final ans. based on all results
eg → by majority (not the best way).

By weighted avg (weight \propto 1/error of that model).

a. Bagging

Consider models

$$M_1, M_2, \dots, M_k$$

Each model should be at least mediocre.

These models are called weak learners. (We choose weak learners to increase diversity).

- * If we give each model a lot of samples, then no. of combos are less so almost same models \Rightarrow no diversity
- If we give too few samples then insufficient to train.
; We usually choose around 70% of total.

- ① To increase diversity, random sampling WITH replacement (called instance/sample bagging).

Weak learner \rightarrow Does better than chance
(eg → for binary, $\frac{1}{2} 50\%$).

- ② Another way to increase diversity is by giving a subset of FEATURES to each model (resulting in diff results).
[feature tagging].

\rightarrow WITHOUT replacement.

If every model is a decision tree, this is called
RANDOM FOREST.

Q. Why does a collection of weak learners do better?
(proof).

Ans: Consider M_1, M_2, \dots, M_{10} .

For misclassification, ^{at least} 11 Models should misclassify.
i.e. $P(11 \text{ models misclassify})$.

Assume, accuracy of each model is 0.7.

$$\Rightarrow \text{error} = 0.3.$$

Assume that the predictions follows a binomial distribution.

$$P(\text{misclassification}) = \underbrace{{}^{10}C_{11}}_{\text{success=misclassify}} \times 0.7^{10} \times 0.3^{11} + {}^{10}C_{10} \times 0.7^9 \times 0.3^{12}$$

$$P[\text{exactly 11 misclassify}] = {}^{10}C_{11} \times 0.7^{10} \times 0.3^{11} \\ = \underline{\underline{0.017}}$$

$$P[12 \text{ misclassify}] = 0.3 \times 10^{-3} = 0.0003.$$

their probability keeps decreasing.

$\therefore P[\text{at least 11 misclassify}]$ is v.v low.

$P[\text{at least 11 are correct}]$ is high.

When error is 0.40 (worse than chance)

$$P(\text{11 misclassify}) = 0.135$$

i.e. ~~not~~ fairly high

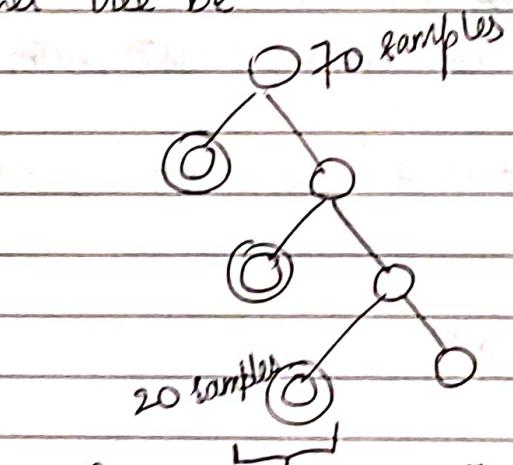
\rightarrow Result has high chance of being wrong.

Q. Why is majority voting not optimal?

Ans: Consider M_1, M_2, M_3, M_4 , all are decision trees.

& $y_i \in \{1, 2, 3, 4\}$. (multi class dataset).

Let tree be



No. Samples: 15, 0, 0, 3
Class 1 2 3 4

\Rightarrow Probabilities: $\frac{15}{20}, \frac{0}{20}, \frac{0}{20}, \frac{3}{20}$

Label distribution

Look at this for all models.

e.g.: 1st 2nd 3rd 4th

	1	4	0	3
M_1	1	4	0	3
M_2	1	4	3	2
M_3	4	3	2	1
M_4	3	4	0	1

If simple majority vote,

we only look at 1st column

& discard 2nd, 3rd & 4th.

2

Output = Class 1

Instead, we should do:

1st col gets $K-1$ votes
 \hookrightarrow no. of classes.

last col gets 0 votes

e.g., here,

1st \rightarrow 3 votes
 2nd \rightarrow 2 votes
 3rd \rightarrow 1 vote
 4th \rightarrow 0 votes

this method
is called
Borda
Count.

Class 1 \rightarrow 6 votes

2 \rightarrow 3 votes

3 \rightarrow 6 votes

4 \rightarrow 9 votes

\therefore Class 4 is output.

* Q. When you can't use Borda Count?

We assume all models are equally strong predictors

So what to do if they aren't equal?

Bagging

\hookrightarrow Instance (w/ replacement)
 \hookrightarrow Feature (w/o replacement).

Instance

It chooses N samples randomly & gives high imp. to
more repeating.

flaw

\hookrightarrow It is chance-based so imp. is arbitrary.

Boosting

We should focus on hard-to-classify samples (hard to learn).

e.g.: let models be

$$M_1, M_2, \dots, M_K$$

w/training set

$$T_1, T_2, \dots, T_k$$

~~Idea~~ Consider 5 samples in dataset:

$$\{x_1, x_2, x_3, x_4, x_5\}$$

Multi-dimensional $\rightarrow d$ dimensions/features

First, build T_i the same way as random forest

(instance bagging)

e.g. \rightarrow take 3 samples.

let $\alpha_i^j \rightarrow$ importance

$P(\text{sample } i \text{ gets chosen for } j^{\text{th}} \text{ tree})$

$$\alpha_1^j = 1/5, \alpha_2^j = \alpha_3^j = \alpha_4^j = \alpha_5^j = 1/5$$

$$\text{e.g. } \rightarrow T_1 = \{x_1, x_2, x_3\}$$

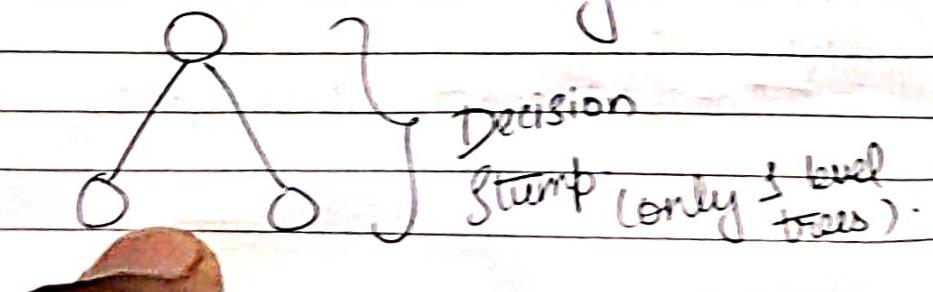
Use this to make a decision tree M_1 .

Now, test M_1 on all 5 samples.

If x_4 is misclassified by M_1 , give it more imp.
In boosting, for decision tree, we ONLY go till ONE level, be:

a) No overfitting

b) No strong learners.



eg: let x_4 be correctly predicted
 & x_3, x_5 be misclassified.

* Every decision step we build should be a weak learner i.e
 $\text{error}(m_j) < 0.5$

Now, we update α_i^2 .

α_3^2 & α_5^2 increase

others decrease.

How to increase/decrease?

[~~Boosting~~ Ada Boost algorithm].

Adaptive

by Freund &
Schapire

formula:

$$\star \alpha_i^{j+1} = \alpha_i^j \times e^{\text{quantity}}$$

look up
AdaBoost mH

If quantity > 0 , weight increases

quantity < 0 , weight decreases.

i.e for ~~correct~~ > 0 , < 0

$e^{+\text{quantity}}$

$e^{-\text{quantity}}$.

$$\boxed{\text{Quantity} = \ln \left(\frac{1 - \text{Error}_j}{\text{Error}_j} \right)}$$

$$\text{eg: } \alpha_3^2 = \frac{1}{5} \times \frac{0.6}{0.4} \quad | \quad \alpha_5^2 = \frac{1}{5} \times \frac{0.6}{0.4} = 0.3$$

$$\alpha_{1,2,4} = \frac{1}{5} \times \frac{0.4}{0.6} = 0.13$$

Now, we want $\sum_{i=1}^5 \alpha_i^j = 1$, but after updating this is not always true.

so we normalise it.

$$\star \boxed{\alpha_i^{j+1} = \frac{\text{Updated } \alpha_i^j}{\sum \alpha_i^{j+1}}}$$



How to take care of probabilities?

Consider 5 buckets -

B_1	B_2	B_3	B_4	B_5
$[0, 0.25)$	$[0.25, 0.5)$	$[0.5, 0.55)$	$[0.55, 0.75)$	$[0.75, 1]$
$\sim 0-1$	$0.25-0.5$	$0.5-0.55$	$0.55-0.75$	$0.75-1$

* Choose a random no. b/w 0 & 1 & check which bucket it falls into. Choose that sample & use it.

* Gradient Boosting Algorithm -

Bias v/s Variance

Variation in output or degree of polynomial.

$$\frac{w^T x - \hat{y}}{\text{variance}}$$

where

Variation = deviations in predictions of model.

Bias = how far you are from the actual value.

(True hypothesis).

low bias \Rightarrow high variance.

\rightarrow low bias \Rightarrow high variance } Trade off.

& vice versa.

Ideally, low bias & low variance, but this is not possible so we settle in the ~~optimal~~ middle.

e.g.: linear regression — high bias.

logistic " for classification — high bias

KNN for 1 neighbour — ~~high~~ bias

↳ smaller the K, lower the bias.

Decision tree w/ large depth — low bias.

↳ higher the depth, lower the bias.

Neural

• Non-linear function

• Many layers

• Many neurons per layer

• Weighted sum of inputs + bias

• Activation function (e.g., sigmoid)

• Backpropagation for training

• Gradient descent for optimization

• Many parameters (weights and biases)

• High computational complexity

• Good at learning complex patterns

• Can be used for classification and regression

• Deep learning variants (e.g., CNN, RNN)

• Still challenging to interpret and explain

• Requires large amounts of data and computation

• Can be prone to overfitting if not properly regularized

• Many hyperparameters to tune (e.g., learning rate, batch size)

• Still an active area of research and development

• Has revolutionized many fields, including computer vision and natural language processing

• Continues to evolve and improve as new techniques and hardware are developed

• Overall, neural networks have become one of the most powerful tools for machine learning and artificial intelligence

• However, they also pose challenges in terms of explainability and ethical concerns around bias and fairness

• Nonetheless, their impact on society is likely to continue to grow and shape the future of technology

• Overall, neural networks represent a significant advancement in the field of machine learning and artificial intelligence

• They have the potential to revolutionize many industries and improve our lives in countless ways

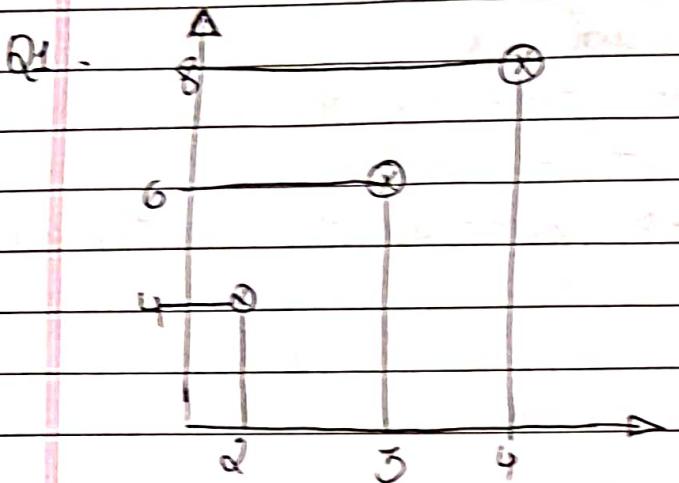
• However, it is important to approach them with a critical eye and consider their limitations and potential risks

• By doing so, we can ensure that they are used ethically and responsibly to benefit everyone

• Overall, neural networks represent a exciting and promising field of study and application

• They have the potential to transform the way we live and work in the years to come

• By continuing to research and develop them, we can help ensure that they are used for the betterment of all

NUMERICALS

x	y
2	4
3	6
4	8

test sample: $x_2 = 5.5$

$$\alpha_0 = \bar{y} - \alpha_1 \bar{x}$$

intuitively, f is αx

$$\alpha_0 = 6 - \alpha(3)$$

$$= 0.$$

$$\Rightarrow y = 2x.$$

(or)

$$(X^T X)^{-1} X^T Y = \begin{bmatrix} \alpha_0 \\ \alpha_1 \end{bmatrix}$$

$$4+6-8$$

$$X = \{x_0, x_1\}$$

$$Y = \{4, 6, 8\}$$

$$\frac{1}{3}$$

$$X = \begin{bmatrix} 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} \quad 3 \times 2$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 \\ 2 & 3 & 4 \end{bmatrix} \quad 2 \times 3$$

$$\frac{1}{54}$$

$$X^T X = \begin{bmatrix} 3 & 9 \\ 9 & 27 \end{bmatrix} = \begin{bmatrix} 3 & 9 \\ 9 & 27 \end{bmatrix}$$

$$\frac{1}{27}$$

$$(X^T X)^{-1} = \begin{bmatrix} 18 & -9 \\ -9 & 6 \end{bmatrix}^{-1} \quad \frac{1}{(54-6)}$$

$$\frac{1}{54}$$

$$(X X^T)^{-1} = \begin{bmatrix} 29 & -9 \\ -9 & -3 \end{bmatrix}^{-1} \times \frac{1}{81-81} = \frac{1}{6}$$

$$\frac{1}{87}$$

$$\begin{bmatrix} 2 & -9 \\ -9 & 3 \end{bmatrix}_{2 \times 2} \times \begin{bmatrix} 1 & 1 & 1 \\ 2 & 3 & 4 \end{bmatrix}_{2 \times 3}$$

$$= \begin{bmatrix} \end{bmatrix}$$

Weight Update Eqⁿ for Gradient Descent

$$\alpha_i^{t+1} = \alpha_i^t - \eta \times \sum_{i=1}^N (\hat{y}_i - y_i) x_i$$

let $\alpha_0 \rightarrow 1$ $\alpha_1 \rightarrow 3$. (assume)] \rightarrow Batch gradient

$$\alpha_0' = \alpha_0^0 - \eta \times \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i) \times x_i$$

take as $\frac{1}{N}$

$$0.01 \quad 3 \quad 1 \quad \hat{y}_i = 1 + 3x$$

$$= 1 - 0.01 \times \frac{1}{3} \times \left[(7-4) + (10-6) + (13-8) \right]$$

$$= 1 - 0.01 \times \frac{1}{3} [12] = 1 - 0.04 = 0.96$$

$$\alpha_1' = 3 - 0.1 \times \frac{1}{3} \times \left[(7-4) + (10-6) + (13-8) \right]$$

≈ 2.88

$$= 3 - 0.1 \times \frac{1}{3} [12] = 3 - 1.2 = 1.74$$

$$= 1.74$$

: keep repeating this

For stochastic, remove $\frac{1}{N}$ & do for 1 eg.

For mini batch, batch size $\frac{N}{M} < N$

Stacking

Take a model on training set.

Assume it predicts \hat{y} .

$$f(x_1, x_2, \dots, x_D) \rightarrow \hat{y}$$

then, train another one on \hat{y}
 $f(\hat{y}) \rightarrow y$

* In d^{rd} model, we can also use it in addition
 to original features x_1, \dots, x_D .

eg: Decision tree, $f(x_1, \dots, x_D) \rightarrow y_{DT}$

Log regression, $f(x_1, \dots, x_D, \hat{y}) \rightarrow y_{LR}$

SVM, $f(y_{DT}, \hat{y}_{LR}) \rightarrow y$

Blending

Instead of simply merging, divide ~~it~~ into partitions
 Each model uses the ~~some~~ data ONLY from 1 partition as
 training sets.

Hierarchical Mixture of Experts

eg: LR $\rightarrow M_1$

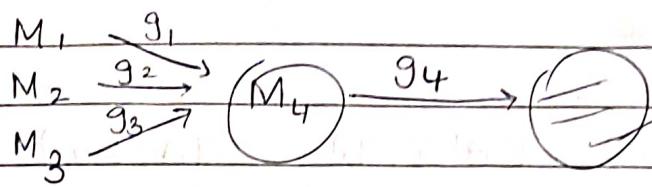
DT $\rightarrow M_2$

SVM $\rightarrow M_3$

$$g_1, g_2, g_3 \rightarrow \sum g_i m_i \quad \text{where } \sum g_i = 1$$

$$\hat{y}_i = \sum g_i m_i$$

$$\min \sum (y_i - \hat{y}_i)^2$$



We can derive g_i^t by gradient descent.

$$g_j^{t+1} = g_j^t - \eta \times \sum_{i=1}^n [\hat{y}_i - y_i] m_j^i$$

prediction

g_i depends on how imp. the model is.

g_i : (if we use mean sq. error),

$$\text{imp. of } M_i = \frac{\sum \text{MSE}_j}{\text{MSE}_{g_i}}$$

$$g_i = \frac{\sum \text{MSE}_j}{\text{MSE}_{g_i}}$$

we need to normalise this so $\sum g_i = 1$.

$$\text{Final weight, } w_i = \frac{g_i}{\sum g_i}$$

$$\therefore w_i = \frac{g_i}{\sum g_i} \quad \text{where } g_i = \frac{\sum \text{MSE}_j}{\text{MSE}_{g_i}}$$

$$\text{& } \hat{y}_i = \sum w_i M_i$$

~~Only for classification~~

Baye's theorem

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

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

Naive Baye's theorem

e.g.

$$P(\text{glucose} = m, \text{age} = 31-40, \text{weight} = H | Y=1) \times P(Y=1) / P(m, 31-40, H)$$

↑ likelihood ↑ prior probability
↓ evidence

$$\text{posterior} = P(Y=1 | m, 31-40, H).$$

But if this sample doesn't exist in training set, $P=0$
To avoid this,

Laplace Correction

Naive :

Give it a v. small probability $\frac{1}{N}$

We assume features are all independent.

so we write

$$P(m, 31-40, H | Y=1) = P(m | Y=1) \times P(31-40 | Y=1) \times P(H | Y=1).$$

[For numeric features, we assume normal distribution]

Consider a numeric feature (eg: BMI),

$$P(M | Y=1) \times P(31-40 | Y=1) \times P(H | Y=1) \times P(35.6 | Y=1)$$

how to find
-1? \hat{x} ?

take mean & std. deviation of training set

use,

$$-(x_i - \mu)^2 / \sigma^2$$

$$\rightarrow P = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

& compute probability
J Gaussian
naive Bayes.

$$\Rightarrow \text{Evidence} = \sum_{j=1}^k P(M, 31-40, H, 35.6 | Y=j) \times P(Y=j).$$

(alternative form)

EVALUATION METRICS

Regression Metrics

$$1) \text{Mean squared error, } MSE = \frac{1}{N} \sum (\hat{y}_i - y_i)^2$$

$$2) \text{Sum of squared errors, } = N \times MSE$$

$$3) \text{Mean absolute error, } \frac{1}{N} \sum |\hat{y}_i - y_i|$$

4) R square (co-efficient of determination). Used for goodness of the fit.

How much the f^n can adapt

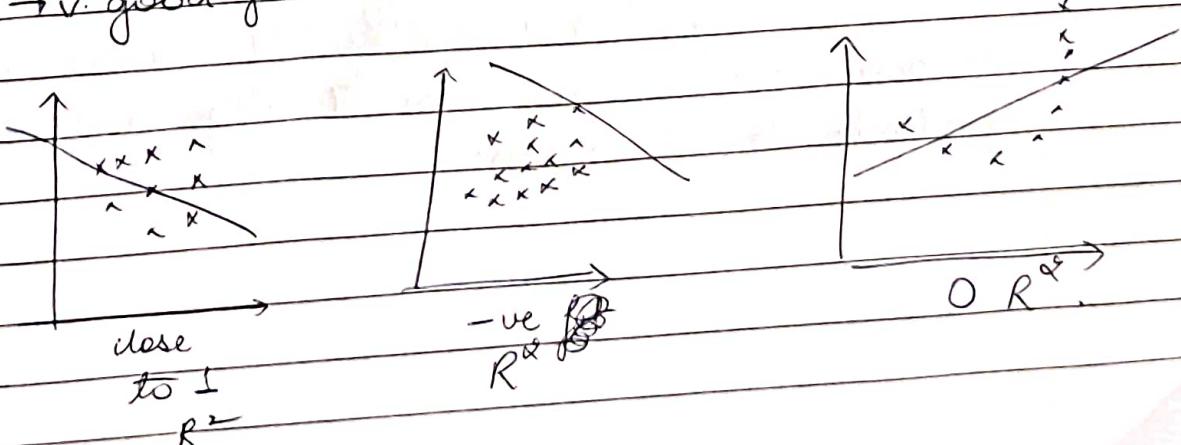
It exists b/w a given range.

It is used to compare diff. degree of polynomials.

-ve → it hasn't learned so far

0 → not a good fit

1 → v. good fit.



$$R^2 \text{ error} = 1 - \frac{\sum_{i=1}^N (\hat{Y}_i - Y_i)^2}{\sum (Y_i - \bar{Y})^2}$$

↓ sum of squared errors/residuals.
 ↓ total sum of squares

If $\hat{Y}_i = \bar{Y}$ i.e. only α_0 i.e. underfitting
 $R^2 = 0$.

If we predict perfectly,
 $R^2 = 1$. (Since predicted = actual)

If sq. error is v.v large
 $R^2 < 0$.

* Only for training set, NOT for test set.

* R^2 is the ONLY metric where higher is better.

Classification Metrics

Indicator variable $[\hat{Y}_i == Y_i]$.

[if equal, 1 else 0].

$$\textcircled{1} \text{ Accuracy} = \frac{\sum [\hat{Y}_i == Y_i]}{N}$$

Error = 1 - accuracy.

We don't prefer accuracy when data is skewed.

Eg.: 99.9% healthy & 0.1% cancerous

We only use accuracy when all categories have almost equal proportion of samples.

(2)

true negative
Predicted

Confusion
Matrix

		0	1	
		TN	FP	→ false positive
True	0	FN	TP	↓ true positive
	1	FP	TN	

$$\text{Here, } (a) \text{accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

~~(b)~~

$$(b) \text{precision} = \frac{TP}{TP+FP}$$

$$(c) \text{recall} = \frac{TP}{TP+FN}$$

In case of skewed data,

try to MAX. ~~BT~~ TP / MIN. FP, based on the situation & requirement.

(3) Binary cross entropy loss (for non-skewed data)

$$= \sum_{i=0}^{\text{samples}} y_i \log o_i - (1-y_i) \log(1-o_i)$$

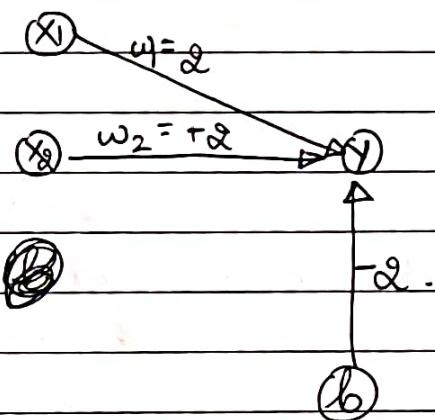
$$\Rightarrow \sum_{\substack{j \text{ over} \\ \text{classes}}} -y_j^j \log o_j$$

(YK-1)
MK

Q. Design AND gate using Hebbian learning.

inputs	x_1, x_2	bias b	y	Outputs			w_1 new	w_2 new	b new
				Δw_1	Δw_2	Δb			
S1	1 1	0	1	1	1	1	1	1	1
	1 -1	0	-1	-1	+1	-1	0	2	10
	-1 1	0	-1	+1	-1	-1	1	+1	-1
	-1 -1	0	-1	1	1	-1	2	2	-2

final weights



Q.

<u>Letters</u>	<u>T₁</u>	<u>T₂</u>	<u>T₃</u>	<u>T₄</u>	<u>T₅</u>	<u>T₆</u>	<u>T₇</u>	<u>T₈</u>	<u>T₉</u>	<u>T₁₀</u>
I	-1	-1	1	-1	1	-1	1	-1	1	0
O	-1	-1	1	-1	1	-1	1	-1	1	0

OBITUARY

December

200
300

1. Surprise 1.22

1

~~11/11 12/12 13/13 14/14 15/15 16/16 17/17 18/18 19/19~~

工 ト ハ ハ ハ ハ ハ ハ ハ ハ

O 21.1.19 - 12.1.20

1. श्रीमद्भागवत

9

Ostb.

وَالْمُؤْمِنُونَ الْمُؤْمِنُونَ الْمُؤْمِنُونَ الْمُؤْمِنُونَ الْمُؤْمِنُونَ

I	-1	+1	-1	+1	-1	-1	-1	-1	-1	1
O	0	0	0	0	2	2	0	0	0	0

100

100₂

$$\alpha \varrho_2 = 0$$

100. -

100

10

10

10

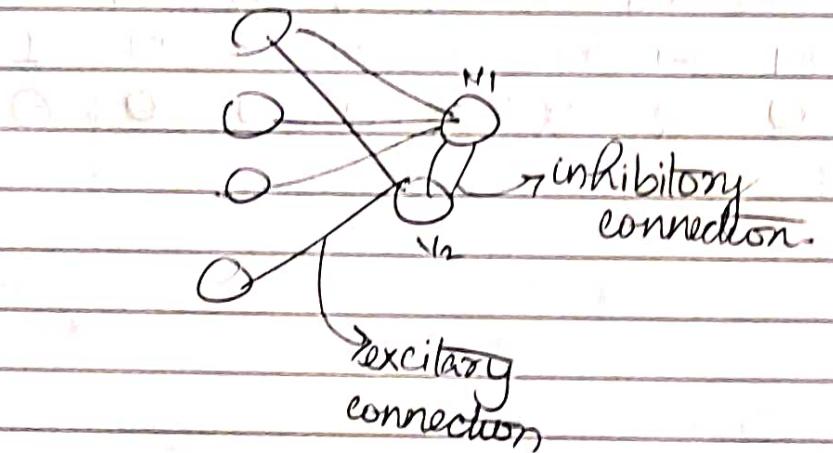


10

6

Competitive Learning

- * Winner-takes-all!
- * Only winner's weights updated



$$v_k = \sum x_j w_{kj}$$

then,

$$\Delta w_{kj} = \begin{cases} \eta(x_j - w_{kj}) & \text{if } k \text{ wins} \\ 0 & \text{if } k \text{ loses} \end{cases}$$

e.g.: $\eta = 0.5$ y_1, y_2

Inputs $\rightarrow [0\ 0\ 1\ 1], [1\ 0\ 0\ 0], [0\ 1\ 1\ 0], [0\ 0\ 0\ 1]$

$$\sum w_{kj} = 1$$

Initialise the weights b/w 0 to 1.

let $w = \begin{bmatrix} 0.2 & 0.9 \\ 0.4 & 0.7 \\ 0.6 & 0.5 \\ 0.8 & 0.3 \end{bmatrix}$

} input layers

} outputs

$$\begin{aligned} v_1 &= [0 \ 0 \ 1 \ 1] \times 0.2 + [1 \ 0 \ 0 \ 0] \times 0.4 + [0 \ 1 \ 1 \ 0] \times 0.6 \\ &\quad + [0 \ 0 \ 0 \ 1] \times 0.8 \end{aligned}$$

$$= [0 \ 0 \ 0.2 \ 0.2] + [0.4 \ 0 \ 0.0] + [0 \ 0.6 \ 0.6 \ 0] + [0 \ 0 \ 0 \ 0.8]$$

$$v_1 = [0.4 \ 0.6 \ 0.2 \ 1.0]$$

Pattern 1

$$v_1 = 0 \times 0.2 + 0 \times 0.4 + 1 \times 0.6 + 1 \times 0.8$$

$$= \underline{\underline{1.4}}$$

v_1 wins

$$\begin{aligned} v_2 &= 0 \times 0.9 + 0 \times 0.7 + 1 \times 0.5 + 1 \times 0.3 \\ &= \underline{\underline{-0.8}} \end{aligned}$$

Pattern 2

$$v_1 = \underline{\underline{0.2}}$$

$$v_2 = \underline{\underline{0.9}}$$

v_2 wins

Pattern 3:

$$v_1 = \underline{\underline{1.0}}$$

$$v_2 = \underline{\underline{1.2}}$$

v_2 wins

Pattern 4

$$v_1 = \underline{\underline{0.8}}$$

$$v_2 = \underline{\underline{0.3}}$$

v_1 wins

PATTERN 1

v_1 wins

$$\Delta w_{11} = 0.5 \times (0 - 0.2) \\ = -\underline{\underline{0.1}}$$

$$\Delta w_{12} = 0.5 \times (0 - 0.4) \\ = -\underline{\underline{0.2}}$$

$$\Delta w_{13} = 0.5 \times (1 - 0.6) \\ = \underline{\underline{0.2}}$$

$$\Delta w_{14} = 0.5 \times (1 - 0.8) \\ = +\underline{\underline{0.1}}$$

$$\rightarrow w_{11} = \underline{\underline{+0.1}} \quad w_{12} = \underline{\underline{+0.2}} \quad w_{13} = \underline{\underline{+0.8}} \quad w_{14} = \underline{\underline{+0.9}}$$

Since node 2 loses, there is no change in weight.

20

Pattern 2V₂ wins

$$\Delta w_{21} = 0.5(1-0.9) \quad \Delta w_{22} = 0.5(1-0.7) \quad \Delta w_{23} = 0.5(1-0.5) \quad \Delta w_{24} = 0.5(1-0.3)$$

$$= 0.05 \quad = -0.35 \quad = -0.25 \quad = -0.15$$

$$w_{21} = 0.95$$

Pattern 3:V₂ wins

$$w_{22} = +0.35$$

$$\Delta w_{21} = 0.5(0-0.9) \quad \Delta w_{22} = 0.5(1-0.7) \quad w_{23} = +0.25$$

$$= -0.45 \quad = 0.15$$

$$w_{24} = +0.15$$

Since node 1 lost, no change.

New weight matrix:

$$W = \begin{bmatrix} 0.1 & 0.95 \\ 0.2 & 0.35 \\ 0.8 & 0.25 \\ 0.9 & 0.15 \end{bmatrix}$$

100
-35
65

Pattern 3:V₂ wins

$$\Delta w_{21} = 0.5(0-0.95) \quad \Delta w_{22} = 0.5(1-0.55) \quad \Delta w_{23} = 0.5(1-0.25) \quad \Delta w_{24} = 0.5(1-0.15)$$

$$= -0.475 \quad = +0.325 \quad = +0.375 \quad = -0.075$$

Pattern 4 V_1 wins

$$\Delta w_1 = 0.5(0 - 0.1) \quad \Delta w_2 = 0.5(0 - 0.2) \quad \Delta w_3 = 0.5(1 - 0.8) \quad \Delta w_{14} = 0.5(1 - 0.9)$$

$$= -0.05 \quad = -0.10 \quad = 0.10 \quad = -0.05$$

$$w_1 = 0.05 \quad w_2 = 0.00 \quad w_{13} = 0.35 \quad w_{14} =$$

Final weight matrix

$$W = \begin{bmatrix} -0.05 & 0.95 \\ -0.10 & 0.35 \\ -0.10 & 0.25 \\ 0.05 & 0.15 \end{bmatrix}$$

New weights after pattern 4

$$W = \begin{bmatrix} 0.1 & 0.95 \\ 0.2 & 0.35 \\ 0.8 & 0.25 \\ 0.9 & 0.15 \end{bmatrix}$$

for pattern 5,

$$V_1 = 0 + 1 \cdot 0.2 + 1 \cdot 0.8 + 0 = 1.0$$

$$V_{02} = 0 + 0.35 + 0.25 + 0 = 0.60$$

 V_1 wins

$$\Delta w_1 = 0.5(0 - 0.1) \quad \Delta w_2 = 0.5(1 - 0.2) \quad \Delta w_3 = 0.5(1 - 0.8) \quad \Delta w_{14} = 0.5(0 - 0.9)$$

$$= -0.05 \quad = 0.40 \quad = 0.10 \quad = -0.45$$

$$W = \begin{bmatrix} 0.05 & 0.95 \\ 0.6 & 0.35 \\ 0.9 & 0.25 \\ 0.45 & 0.15 \end{bmatrix}$$

Pattern 4

$$v_1 = 0.45 \quad v_2 = 0.15$$

 v_1 wins

$$\Delta w_{11} = 0.5(-0.05) \quad \Delta w_{12} = 0.5(-0.6) \quad \Delta w_{13} = 0.5(-0.9) \quad \Delta w_{14} = 0.5(0.55)$$

$$= -0.025 \quad \quad \quad = -0.3 \quad \quad \quad = -0.45 \quad \quad \quad = 0.275$$

$$W = \begin{bmatrix} 0.025 & 0.95 \\ 0.3 & 0.35 \\ 0.45 & 0.25 \\ 0.725 & 0.15 \end{bmatrix}$$