

## \* K-Nearest neighbour.

- Also known as lazy learner/ instance-based learning.
- Used when data is distributed in non-linear manner.

$K \Rightarrow$  how many nearest values we need to consider in terms of distance

↳ Euclidean  
↳ Manhattan

### ① for classification use-case

- The answer for the new point would be the one with majority of category. and thus gets classified to that group

### ② Regression use-case.

- Here all values will be of same group
- We will find the mean of all  $k$  nearest neighbours.

- KNN will get impacted by outliers, because if there are a lot of outliers, then it will be difficult for a classification prediction.
- Same will happen if the data is imbalanced. and leads to bias prediction.

- |  |   |
|--|---|
| <ul style="list-style-type: none"> <li>→ static</li> <li>→ RIP</li> <li>→ OSPF</li> <li>→ subnetting</li> <li>→ virtual LAN</li> <li>→ wireless LAN</li> </ul> | <ul style="list-style-type: none"> <li>→ When dataset is too big, classification time is too high.</li> <li>→ When <math>K</math> is high, the less similar data (one whose value is 0 if done with cosine similarity) will contribute which leads to wrong prediction.</li> <li>→ When <math>K</math> is too small, wrong prediction may lead as if one outlier of one class is in other coordinate whose other class has so many points, this causes wrong prediction.</li> </ul> |
|--|---|

SE =

ML - 14

DAA -

- $k$  is small  $\Rightarrow$  overfit
- $k$  is large  $\Rightarrow$  underfit
- **Thumb rule**  $\Rightarrow k = \sqrt{n}$
- **Note:**  $n$  is no. of training examples.
- \*  $k$  should be odd.
- \*  $k$  value must not be multiple of no. of classes.
- \* Should not be too small or too large
  - too small  $\Rightarrow$  greater influence to outlier/noise.
  - too big  $\Rightarrow$  bias towards highly probable class.
- In KNN classifier, we take avg. of  $k$  near values.

- \* **Naive Bayes classifier.**
- Used bayes theorem, i.e., probabilistic approach to clas the data.
- They are generally used for text classification or medical image related classification.
- A fast machine learning model for quick prediction.

- Two assumptions.
  - features given are independent of each other.
  - All the features contribute equally.

$$\text{Posterior probability } P(C|x) = \frac{\text{likelihood}}{\text{prior probability}} \cdot P(x|C) \cdot P(C) \rightarrow \text{prior probability}$$

$C = \text{Class}$

$x = \text{feature/test sample}$

→ Here  $P(x)$  is constant ( $=1$ ) for all classes

→ Let  $D$  be no. of training set of tuples and has  $n$  features  
∴ we will form  $n$ -D attribute vector  $X = (x_1, x_2, \dots, x_n)$

→ Suppose there are  $m$  classes  $C_1, C_2, \dots, C_m$   
Classification is to derive max. posteriority; i.e., The max.  $P(C_i|x)$   
∴  $P(C_i|x) = \frac{P(x|C_i)P(C_i)}{\sum P(x)}$

BUT  $P(x) = \text{constant} = 1$   
 $P(C_i|x) = \underbrace{P(x|C_i) * P(C_i)}$   
needs to be maximized

Now,  $P(x|C_i) = \prod_{k=1}^n P(x_k|C_i) = P(x_1|C_i) \cdot P(x_2|C_i) \cdot P(x_3|C_i) \cdots P(x_n|C_i)$

→ If  $A_k$ /feature( $k$ ) is categorical,  $P(x_k|C_i)$  is the no. of tuples in  $C_i$  having value  $x_k$  for  $A_k$  divided by  $|C_i|$ .

→ If  $A_k$  is continuous value,  $P(x_k|C_i)$  is computed based on gaussian distribution with mean  $\mu$  and S.D.  $\sigma$ .

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

and  $P(x_k|C_i)$  is  $= g(x_k, \mu_{C_i}, \sigma_{C_i})$ .

# Avoiding zero probability problem.

→ We use laplacian correction.

Suppose for 1000 tuple, income = low(1), medium(99), high(10)

∴ we add 1 to each case

So, now, tuple income = low(1), medium(99), high(11)  
and no. of tuple will be 1003.

## \* Weighted KNN

- first find the k-nearest neighbor based on distance.
- Then, inverse them.
- Then, sum the ~~inverse~~ inverted (done in step 2) K neighbours.
- Then, redefined the weights ~~vector~~, i.e.,  $\text{inversed}(\text{step 2}) / \text{sum}$   
for each of the k-neighbors.
- Then sum them based on their classes.  
like class 0 =  $\text{sum}_1$ ; class 1 =  $\text{sum}_2$  . . .
- Then one with highest sum, the test data is assigned that class.

## \* Model Evaluation and Selection.

- Methods for estimating a classifier's accuracy.

- (a) Holdout method, random subsampling
- (b) cross-validation
- (c) Bootstrap.

- Validation techniques.

### \* Resubstituting

#### \* Hold out

- LOOCV (Leave-One-Out-Cross-Validation)

- Random subsampling

- Bootstrapping

\* Classifier Evaluation Metrics = Confusion Matrix.

AC   PC	$C_i$	$\neg C_i$	
$C_i$	TP	FN	P
$\neg C_i$	FP	TN	N

$(CM_{ij} \Rightarrow)$  no. of tuples in class i that were labelled by the classifier as class j.

$$\rightarrow \text{Accuracy} = (TP + TN) / AU$$

$$\rightarrow \text{Error rate} = 1 - (\text{accuracy}) = (FN + FP) / AU$$

$$\rightarrow \text{Sensitivity} = \text{True positive recognition rate} = TP / P$$

$$\rightarrow \text{Specificity} = \text{True negative recognition rate} = TN / N$$

$$\rightarrow \text{Precision} = \text{what \% of tuple classifier labelled +ve are actually positive}$$

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

$$\rightarrow \text{Recall} = \text{what \% of +ve tuples did the classifier label as +ve?}$$

$$= \frac{TP}{TP + FN}$$

$$\rightarrow \text{Note} = \text{Inverse relation b/w precision \& recall}$$

$$\rightarrow f\text{-measure (F1 score or F score)} = \text{Harmonic mean of precision \& recall.}$$

$$F = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

$$\rightarrow F_\beta : \text{weighted measure of precision and recall.}$$

$$F_\beta = \frac{(1 + \beta^2) * \text{precision} * \text{recall}}{\beta^2 * \text{precision} + \text{recall}}$$

\* Naive Bayes Classifier = Multinomial.

- Compulsory to do Laplacian correction.
- Generally used for document classification, where we are interested in finding the frequency of occurrence of that specific term in that document.

$$P(d_j | C)$$

↓

Probability of  
document belonging to  
class C.

i, m = terms

J, n = documents

$$P(d_j | C) = \frac{\left(\sum_{i=1}^m n_{ij}\right)!}{\prod_{i=1}^m n_{ij}!} \prod_{i=1}^m P(t_i | C)^{n_{ij}}$$

length of document j      no. of appearance of term i  
in document j.

$$P(t_i | C) = \frac{\sum_{j=1}^n n_{ij}}{\sum_{i=1}^m \sum_{j=1}^n n_{ij}}$$

(i) no. of appearance of ith term in all  
documents of class C.

← sum of length of all documents of class C.

→ Resubstitution error

When we are using the set on which we train our model and we are testing the same set again, for accuracy. Then the error we get is called resubstitution error.

→ Bootstrapping (0.632 bootstrap)

from ppt.

## \* Decision Tree.

→ Decision tree induction is the learning of decision trees from class-labeled training tuples.

→ \* Each internal node denotes the test on an attribute.

\* Each branch/edge represents outcome of test.

\* Each leaf node holds a class label.

→ DT has good classification accuracy.

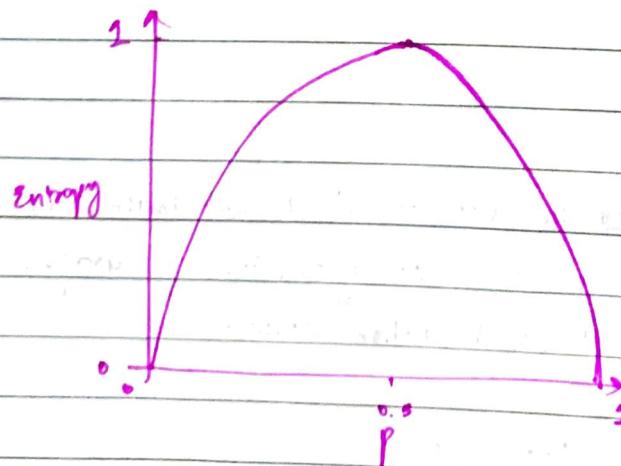
→ \* When outcomes are certain, we will get less information.

\* When outcomes are uncertain, we will get information.

→ Entropy is the measure of impurity of dataset D.

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2(p_i)$$

\* As the data becomes purer and purer, the entropy value becomes smaller and smaller.



## \* Classification by Decision Tree Induction - ID3.

→ See lecture of 20/09/2021

- The problem with ID3 classification is that it is biased towards tests with many outcomes, i.e., it prefers to select attributes having a large no. of values.
- Thus this will help in classification better, but such a partitioning used could be useless.

## \* Classification by Decision Tree Induction - C4.5.

- C.4.5 is successor of ID3 which overcomes this bias of ID3, using gain ratio.

### → Gain ratio

- Kind of normalization to information gain using a split information value defined analogously with  $\text{Info}(D)$  as

$$\text{SplitInfo}_A(D) = \sum_{j=1}^k \frac{|D_{ij}|}{|D|} \log_2 \left( \frac{|D_{ij}|}{|D|} \right)$$

↑ no. of attributes.

$$\text{Gain Ratio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}(A)}$$

- The attribute with max. gain ratio is selected as splitting attribute.

g

WB	F	fur	S	Lays Egg
Y	Y	N	No	Y
N	N	N	Y	Y
Y	Y	N	N	Y
Y	Y	N	N	Y
Y	N	N	Y	N
Y	N	Y	N	N

Using DD 3

$$\text{Info}(D) = -\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right)$$

$$= 0.9182 \text{ bits}$$

$$2\text{Info}_{WB}(D) = \frac{5}{6} \left[ -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \right] + \frac{1}{6} \left[ -1 \log_2\left(\frac{1}{1}\right) \right]$$

$$= 0.8091 \text{ bits}$$

$$\text{Info}_F(D) = \frac{3}{6} \left[ \frac{3}{3} \log_2\left(\frac{3}{3}\right) \right] + \frac{3}{6} \left[ -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) \right]$$

$$= 0.4591 \text{ bits}$$

$$\text{Info}_{\text{fur}}(D) = \frac{1}{6} [D] + \frac{5}{6} \left[ -\frac{1}{5} \log_2\left(\frac{1}{5}\right) - \frac{4}{5} \log_2\left(\frac{4}{5}\right) \right]$$

$$= 0.6016 \text{ bits}$$

$$\text{Info}_S(D) = \frac{2}{6} \left[ -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \right] + \frac{4}{6} \left[ -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right]$$

$$= 0.8741 \text{ bits}$$

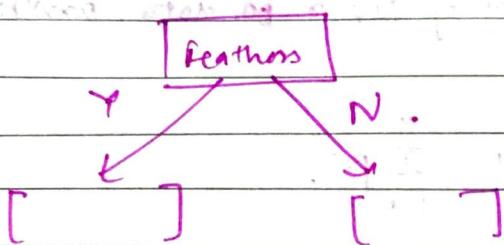
$$\therefore \text{Gain}(WB) = 0.9182 - 0.8091 = 0.109 \text{ bits}$$

$$\text{Gain}(F) = 0.9182 - 0.4591 = 0.4591 \text{ bits}$$

$$\text{Gain}(x_{\text{min}}) = 0.9152 - 0.6016 \approx 0.3164$$

$$\text{Gain}(s) = 0.9152 - 0.8741 = 0.0441.$$

∴ Max gain is of  $f$ .



## \* Computing Information-Gain for continuous-valued attributes.

→ Let attribute A be continuous-valued attribute.

\* discretization (break attribute into ranges in advance)

\* using thresholds for splitting nodes.

→ Select point with min. expected information required.

Length	10	15	21	28	32	40	50
Class	N	Y	Y	N	Y	Y	N

$$A_1 = \frac{10+15}{2}$$

$$\text{no. of samples in partition} = 12.5$$

$$A_2 = \frac{21+28}{2}$$

$$= 24.5$$

$$A_3 = \frac{32+40}{2}$$

$$= 36$$

$$\text{Info}_{A_1}(D) = \frac{1}{7} \left[ \frac{-1}{1} \log \frac{1}{1} \right] + \frac{6}{7} \left[ \frac{-4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6} \right]$$

no. in partition 1    no. in partition 2

$$\text{Info}_{A_2}(D) = \frac{3}{7} \left[ \frac{-2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} \right] + \frac{4}{7} \left[ \frac{-2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} \right]$$

$$\text{Info}_{A_3}(D) = \dots$$

∴ Select point with min. info.

## \* Classification by Decision Tree Induction - CART

→ Gini index.

- Used when we want to have binary partition
- measure of impurity of  $D$ , a  $\text{pa}$  data partition or set of training tuples. as

$$\text{Gini}(D) = 1 - \sum_{i=1}^m p_i^2.$$

→ for discrete value attribute, all possible combination of value except empty set & power set.

$$\text{Gini}_A(D) = \frac{|D_1|}{D} \text{Gini}(D_1) + \frac{|D_2|}{D} \text{Gini}(D_2) - \dots$$

→ subset that gives min. gini index is selected as splitting subset.

→ for continuous, each possible split point must be considered.

→ Reduction in impurity that would be incurred by a binary split on a discrete/continuous attribute  $A$  is,

$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D).$$

→ Attribute that max. the reduction impurity is selected as splitting attribute.

→ Information gained is biased towards multivalued attribute.

→ Pruning of decision tree is used to avoid overfitting.  
Overfitting when no. of features are high.

## \* Support Vector Machine (SVM)

- Used for binary classification.
- Used for high dimensional data.
- SVM is linear learning system that builds two class classifier.

$$\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)\} \Rightarrow D_{\text{Training}}$$

$$x_i \Rightarrow (x_{i1}, x_{i2}, \dots, x_{ir}) \quad x \in \mathbb{R}^r$$

$\curvearrowleft$  this are the features

$y_i$  is a class label  $y_i \in \{-1, 1\}$

SVM finds linear function which is of the form

$$f(x) = \langle w, x \rangle + b$$

If  $f(x_i) \geq 0$ , it results in class label 1.  $\therefore y_i = 1$   
 If  $f(x_i) < 0$ , then class label is -1.  $\therefore y_i = -1$ .

$$\therefore y_i = \begin{cases} 1, & \text{if } f(x_i) \geq 0 \\ -1, & \text{if } f(x_i) < 0 \end{cases}$$

$$f(x_1, x_2, \dots, x_r) = \underbrace{w_1 x_1 + w_2 x_2 + \dots + w_r x_r}_b + b$$

this is a hyperplane.

hyperplane in 2D  $\Rightarrow$  line

hyperplane in 3D  $\Rightarrow$  plane.

- Disadvantage of having min distance, margin is that we want to avoid false prediction

- We select a hyperplane which is away from both the classes equally.

A large margin effectively corresponds to a regularization of SVM weights which prevents overfitting. Hence, we prefer a large margin (or the right margin chosen by cross-validation) because it helps us generalize our predictions and perform better on the test data by not overfitting the model to the training data.

$$\rightarrow w \cdot x + b = 0$$

$w \rightarrow$  defines a direction  $\perp$  to the hyperplane. (normal vector).

→ Maximum margin hyperplane

↳ It gives you maximum margin b/w the +ve and -ve class samples.

$$H_+ = \langle w \cdot x^+ \rangle + b = 1$$

$$H_- = \langle w \cdot x^- \rangle + b = -1$$

such that  $\langle w \cdot x_i \rangle + b \geq 1$  if  $y_i = 1$

$\langle w \cdot x_i \rangle + b \leq -1$  if  $y_i = -1$

$$(d_+) + (d_-) = \text{max. margin.}$$

$$\text{and } (d_+) + (d_-) = \frac{2}{\|w\|}$$

We reverse this,

$$\text{thus, minimize } \frac{\langle w \cdot w \rangle}{2} / \frac{\|w\|^2}{2}$$

$$\text{and } y_i(\langle w \cdot x_i \rangle + b) \geq 1 \text{ for all } i$$

⇒ Standard Lagrangian Multiplier Method.

$$\text{minimize } \frac{\langle w \cdot w \rangle}{2}$$

$$\text{subject to } y_i(\langle w \cdot x_i \rangle + b) \geq 1, i = 1, 2, \dots, n$$

same as  
this equation

$$\therefore L_p = \frac{1}{2} \langle w, w \rangle - \sum_{i=1}^n \alpha_i [y_i (\langle w, x_i \rangle + b) - 1]$$

Constraints  $w_1, w_2, w_3, \dots, w_n$  and  $b$ .

Now,  $\frac{dL_p}{dw_j} = w_j - \sum_{i=1}^n y_i \alpha_i x_{ij} = 0$  where  $j = 1, 2, \dots, n$

$$\frac{dL_p}{db} = - \sum_{i=1}^n y_i \alpha_i = 0.$$

$$(y_i (\langle w, x_i \rangle + b) - 1) \geq 0 \text{ for } i = 1, \dots, n.$$

$$\alpha_i \geq 0 \text{ for } i = 1, 2, \dots, n$$

$$\underbrace{x_i(y_i(\langle w, x_i \rangle + b) - 1)}_0 = 0 \text{ for } i = 1, \dots, n$$

only those points which are on marginal planes for which  $\alpha_i > 0$ .

$\therefore$  those points (samples) are called support vectors.

$\rightarrow \therefore$  for SVM, not all points / data is used for identifying the hyperparameter.

see lectures for better understanding.

Kuch samej nahi data SVM me.

$$L_D = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j x_{ij} x_{ij} \alpha_i \alpha_j$$

wolfe dual

$$\text{subject to } \sum_{i=1}^n y_i \alpha_i = 0$$

$$\alpha_i \geq 0 ; i = 1, \dots, n$$

# SVM

# Support Vector Machine

Eg:

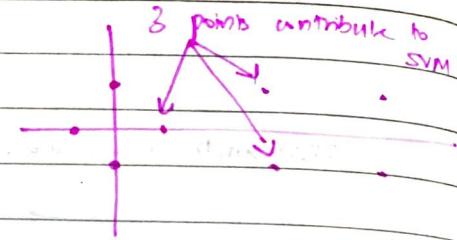
SRM

the labelled points in  $\mathbb{R}^2$ :

$$\{(3,1), (3,-1), (6,1), (6,-1)\}$$

-ve labelled points in  $\mathbb{R}^2$ :

$$\{(1,0), (0,1), (-1,0), (0,-1)\}$$



$$\therefore \vec{x}_1 = (1,0) \text{ by Class -1.}$$

$$\vec{x}_2 = (3,1)$$

$$\vec{x}_3 = (3,-1) \quad \left. \right\} \text{Class 1}$$

The vector is augmented with 1 as "bias input".

$$\therefore \vec{x}_1 = (1,0,1)$$

$$\vec{x}_2 = (3,1,1)$$

$$\vec{x}_3 = (3,-1,1)$$

Objective function.

$$w(x) = \sum_{i=1}^3 \alpha_i \vec{x}_i - \frac{1}{2} \sum_{i=1}^3 \sum_{j=1}^3 \alpha_i \alpha_j y_i y_j K(\vec{x}_i, \vec{x}_j)$$

Kernel function is identity.

$$K(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$

and  $d = 3 = \text{no. of points contribute to SVM}$ .

$$\therefore w(x) = \sum_{i=1}^3 \alpha_i - \frac{1}{2} \sum_{i=1}^3 \sum_{j=1}^3 \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j).$$

$$= \sum_{i=1}^3 \alpha_i - \frac{1}{2} \left[ \alpha_1 x_1 y_1 (\vec{x}_1 \cdot \vec{x}_1) + \alpha_2 x_2 y_2 (\vec{x}_2 \cdot \vec{x}_2) + \dots + \alpha_3 x_3 y_3 (\vec{x}_3 \cdot \vec{x}_3) \right].$$

$$(\vec{x}_1 \cdot \vec{x}) = [1 \ 0 \ 1] \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = 2$$

$$\vec{x}_1 \cdot \vec{x}_2 = [1 \ 0 \ 1] \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} = 4.$$

$$(\vec{x}_1 \cdot \vec{x}_3) = 4; (\vec{x}_2 \cdot \vec{x}_1) = 1; (\vec{x}_2 \cdot \vec{x}_3) = 9; (\vec{x}_3 \cdot \vec{x}_1) = 11.$$

$$\therefore W(x) = \sum_{i=1}^3 \alpha_i - \frac{1}{2} \left[ 2x_1^2 - 4x_1x_2 - 4x_1x_3 - 4\alpha_1\alpha_2 + 11x_2^2 + 9x_2x_3 - 4\alpha_1\alpha_3 + 9\alpha_2\alpha_3 + 11\alpha_3^2 \right].$$

$$= \alpha_1 + \alpha_2 + \alpha_3 - \frac{1}{2} \left[ 2x_1^2 - 8x_1\alpha_2 - 8\alpha_1\alpha_3 + 11x_2^2 + 18x_2x_3 + 11\alpha_3^2 \right].$$

$\therefore$  we need to max. this equation w.r.t  $x$ .

$$\therefore \frac{dW(x)}{dx_1} = 1 - \frac{1}{2} [4\alpha_1 - 8\alpha_2 - 8\alpha_3] = 0.$$

$$2\alpha_1 - 4\alpha_2 - 4\alpha_3 = 1.$$

$$\frac{d(W(x))}{dx_2} = 4\alpha_1 - 11\alpha_2 + 9\alpha_3 = -1$$

$$\frac{d(W(x))}{dx_3} = 4\alpha_1 - 9\alpha_2 - 11\alpha_3 = -1.$$

Solving the above 3 equation

$$\therefore \alpha_1 = -3.5, \alpha_2 = 0.75, \alpha_3 = 0.75.$$

Now we have  $\alpha_i$ , we find hyperplane from

$$\vec{w} = \sum_{i=1}^3 \alpha_i y_i \vec{x}_i$$

$$= 3.5 \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} + 0.75 \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix} + 0.75 \begin{bmatrix} 3 \\ 1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 0 \\ -2 \end{bmatrix}$$

since the vectors are augmented with bias, the last entry in  $\vec{w}$  is value of bias  $b$ , i.e., ~~b = -2~~  $b = -2$  and  $\vec{w} = [1 \ 0 \ -2]$

Note hyperplane is

$$\vec{w} \cdot \vec{x} + b = 0.$$

Since dataset is 2D,

$$\text{then } w_1x_1 + w_2x_2 + b = 0.$$

$$1 \cdot x_1 + 0 \cdot x_2 - 2 = 0$$

$$\boxed{x_1 = 2}$$

Now, if we want to classify  $x = (4, 1)$  then.

$$f(x) = \text{sgn}(\vec{w} \cdot \vec{x} + b)$$

$$f\left(\begin{bmatrix} 4 \\ 1 \end{bmatrix}\right) = \text{sgn}\left((4 \cdot 1)\begin{bmatrix} 1 \\ 0 \end{bmatrix} - 2\right) = \text{sgn}(4 - 2) \\ = \text{sgn}(2)$$

Eg - 2  
~~shortest method~~

$$S_1 = \begin{pmatrix} ? \\ 1 \end{pmatrix}$$

$\hookrightarrow$  Class = -1

$$S_2 = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$$

$\hookrightarrow$  Class = -1

$$S_3 = \begin{pmatrix} 4 \\ 0 \end{pmatrix}$$

$\hookrightarrow$  Class = 1.

$$\alpha_1 S_1 S_1 + \alpha_2 S_2 S_1 + \alpha_3 S_3 S_1 = -1$$

$$\alpha_1 S_1 S_2 + \alpha_2 S_2 S_2 + \alpha_3 S_3 S_2 = -1$$

$$\alpha_1 S_1 S_3 + \alpha_2 S_2 S_3 + \alpha_3 S_3 S_3 = +1$$

} 3 equations

when differentiated

wrt  $\alpha_i$

$$6x_1 + 4x_2 + 9x_3 = -1$$

$$4x_1 + 6x_2 + 9x_3 = -1$$

$$9x_1 + 9x_2 + 17x_3 = 1$$

$$\alpha_1 = -3.25, \alpha_2 = -3.25, \alpha_3 = 3.5$$

$$\therefore \vec{w} = \sum_{i=1}^3 \alpha_i y_i \vec{x}_i$$

$$\vec{w} = -3.25 \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix} - 3.25 \begin{bmatrix} 2 \\ -1 \\ 1 \end{bmatrix} + 3.5 \begin{bmatrix} 4 \\ 0 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 0 \\ -3 \end{bmatrix}$$

$$\therefore \vec{w} = \begin{bmatrix} 1 \\ 0 \\ -3 \end{bmatrix}, b = -3.$$

$\therefore$  Hyperplane

$$\vec{w} \cdot \vec{x} + b = 0$$

Since dataset is 2D.

$$\therefore w_1x_1 + w_2x_2 + b = 0.$$

$$1 \cdot x_1 + 0 \cdot x_2 - 3 = 0.$$

$$\boxed{x_1 = 3}$$

## \* SVM - Inseparable case.

Due to noise/outlier, we need to modify the equation. Thus we add some error term.

$$\therefore \langle \vec{w}, \vec{x}_i \rangle + b \geq 1 - \varepsilon_i \quad \text{for } y_i = 1$$

$$\langle \vec{w}, \vec{x}_i \rangle + b \leq -1 + \varepsilon_i \quad \text{for } y_i = -1.$$

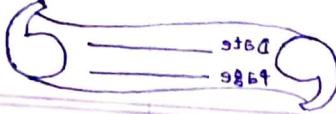
$$\therefore y_i(\langle \vec{w}, \vec{x}_i \rangle + b) \geq 1 - \varepsilon_i \quad \text{for } i = 1, 2, \dots, n$$

$$\varepsilon_i \geq 0, \quad i = 1, 2, \dots, n.$$

$$\text{Also, minimize} = \frac{\langle \vec{w}, \vec{w} \rangle}{2} + C \left( \sum_{i=1}^n \varepsilon_i \right)^k$$

where  $C \geq 0$  is user specified parameter.

Here we will discuss  $k=1$  case only.



Now equation becomes,

$$L_P = \frac{1}{2} \langle w, w \rangle + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i (\langle w, x_i \rangle + b) - 1 + \xi_i] - \sum_{i=1}^n \mu_i \xi_i$$

where  $\alpha_i, \mu_i \geq 0$  are Lagrange multipliers.

And, we can find value of  $b$  as.

$$b = \frac{1}{\sum y_i} - \sum_{i=1}^n y_i \alpha_i \langle x_i, x_j \rangle$$

Again, due to numerical errors, we can compute all possible  $b$ 's and take average  $b$  as final value.

### \* Non-Linear SVM.

→ The basic idea is to map data in the input space  $X$ , to a feature space  $F$  via a nonlinear mapping  $\phi$ .

$$\phi: X \rightarrow F$$

$$x \mapsto \phi(x).$$

After mapping, original dataset  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  becomes  $\{(\phi(x_1), y_1), (\phi(x_2), y_2), \dots, (\phi(x_n), y_n)\}$

### \* Perceptron model.

$$a = \text{hardlim}(w_p + b)$$

$\therefore$  ~~if~~  $a = 1$  if net  $(w_p + b) \geq 0$ ,  
else  $a = 0$  if ~~net~~ " $\leq 0$ "

- $w_{new} = w_{old} + \eta e P \rightarrow$  input samples  
 ↓ learning rate.

$$b_{new} = b_{old} + \eta e, \text{ where } e = \text{target} - \text{actual}$$

Eg: input {0.0, 0.17, 0.33, 0.50, 0.67, 0.83, 1.0} Y

first four are of class 0 (desired output is 0).

Remaining are of class 1 (desired output is 1).

Learning rate = 0.1

Initial  $w = -0.26$

Initial  $b = -0.1$ .

→ The value of the bias must be such that it guarantees it does not cross the origin.

→ Issues

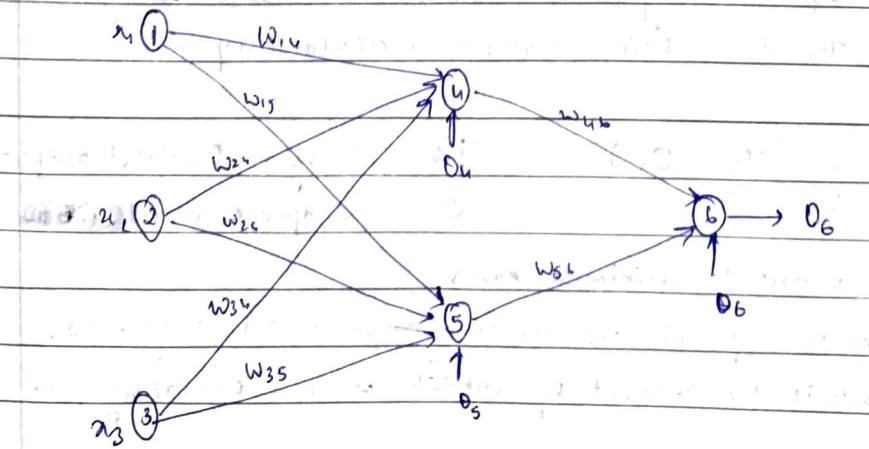
→ Terminating condition  $\Leftrightarrow$  What is the tolerable amount of error.

→ Learning rate

→ Non-numeric inputs

→ Epoch.

### Back propagation



Above is example of multilayer feed forward neural network.

Initial values

$x_1$	$x_2$	$x_3$	$w_{14}$	$w_{15}$	$w_{24}$	$w_{25}$	$w_{34}$	$w_{35}$	$w_{46}$	$w_{56}$	$\theta_4$	$\theta_5$
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2

Class label = 1.

$net(4) =$  input given to neuron 4.

$$= w_{14}x_1 + w_{24}x_2 + w_{34}x_3 + \theta_4$$

$$net(5) = w_{15}x_1 + w_{25}x_2 + w_{35}x_3 + \theta_5.$$

①  $O(4) =$  Output by neuron 4.

$$= f(net(4))$$

$$= \frac{1}{1 + e^{-net(4)}} \quad (\text{Sigmoid function})$$

$$net(6) = w_{46}\theta_4 + w_{56}\theta_5 + \theta_6$$

$$= (-0.3)f(net(4)) + (-0.2)f(net(5)) + 0.1.$$

Backpropagation algorithm, which is generalization of the least mean squared algo. that modifies network weights to minimize the MSE b/w desired & actual output.

$$\text{Error} = \frac{1}{2} (T_d - O_d)^2$$

$T_d$  = Target / Desired output.

$O_d$  = Output at ~~neuron~~ 6.

$\therefore$  we want to minimise error.

$\therefore$  we take differentiation of this error w.r.t weight  $w$ .

and this is the amount by which we are updating the weight.

∴ update in network parameters.

$$w_{44} \quad w_{45} \quad w_{46}$$

$$w_{54} \quad w_{55} \quad w_{56}$$

$$w_{34} \quad w_{35} \quad \Theta_4$$

$$w_{36} \quad \Theta_5 \quad \Theta_6$$

$$\Theta_6$$

then,  $w_{46}, w_{56}, \Theta_4, \Theta_5, \Theta_6$

∴ In backpropagation, we update  $\Theta_6$ , then propagate it backwards.

$$\Theta_6 = \Theta_6 - \eta \frac{dE}{d\Theta_6}$$

$$\text{Now, } \frac{dE}{d\Theta_6} = \frac{dE}{dO_6} \times \frac{dO_6}{d\text{net}_6} \times \frac{d\text{net}_6}{d\Theta_6}$$

$$\frac{dE}{dO_6} = \left[ \frac{1}{2} \times 2 \times (T_6 - O_6)(-1) \right] = -(T_6 - O_6) = O_6 - T_6$$

$$\text{Now } O_6 = f(\text{net}(O))$$

↳ sigmoid.

$$\frac{d f(\text{net}(O))}{d \text{net}(O)} = f(\text{net}(O)) \times (1 - f(\text{net}(O)))$$

$$\frac{d O_6}{d \text{net}_6} = O_6(1 - O_6)$$

$$\text{and } \text{net}(O) = w_{46}O_4 + w_{56}O_5 + \Theta_6.$$

$$\frac{d(\text{net}(O))}{d\Theta_6} = 1.$$

$$\therefore \frac{dE}{d\Theta_6} = (O_6 - T_6) \times O_6(1 - O_6) \times 1 = -(T_6 - O_6) O_6(1 - O_6).$$

Now, for  $w_{46}$ .

$$\frac{dE}{dw_{46}} = \frac{dE}{dO_6} \times \frac{dO_6}{d\text{net}_6} \times \frac{d\text{net}_6}{d w_{46}}$$

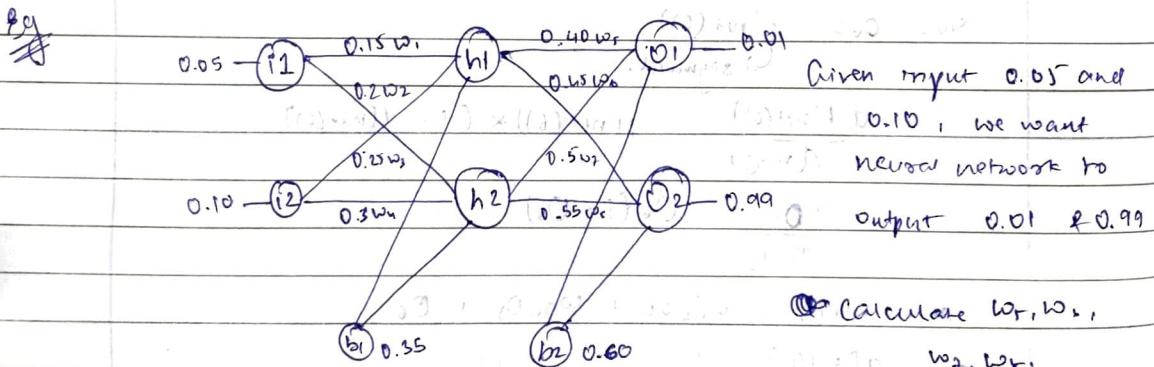
$$= (-1)(T_6 - O_6) O_6(1 - O_6) \times (O_4)$$

$$\frac{dE}{d\omega_{66}} = (-1)(T_6 - O_6) O_6 (1 - O_6) \times O_5.$$

$$\begin{aligned} \frac{dE}{dO_4} &\rightarrow \frac{dE}{dO_6} \times \frac{dO_6}{d\omega_{66}} \times \cancel{\frac{dnet^6}{dO_4}} \times \frac{dnet^6}{dO_4} \times \frac{dO_4}{dnet^4} \times \cancel{\frac{dnet^4}{dO_4}} \\ &= -(T_6 - O_6) O_6 (1 - O_6) \times \omega_{46} \times O_4 (1 - O_4) \times 1 \end{aligned}$$

$$\frac{dE}{d\omega_{46}} = -(T_6 - O_6) O_6 (1 - O_6) \times \omega_{46} \times O_4 (1 - O_4) \times 1.$$

$$\frac{dE}{d\omega_5} = -(T_6 - O_6) O_6 (1 - O_6) \times \omega_{56} \times O_5 (1 - O_5) \times 1.$$



$$\begin{aligned} \text{net } h_1 &= i_1 * w_1 + i_2 * w_2 + b_1 \\ &= (0.05 * 0.15) + (0.1 * 0.20) + 0.35 \\ &= 0.3775 \end{aligned}$$

$$\begin{aligned} \text{net } h_2 &= i_1 * w_3 + i_2 * w_4 + b_2 \\ &= (0.05 * 0.25) + (0.1 * 0.3) + 0.35 \\ &= 0.393 \end{aligned}$$

$$\text{out } o_1 = \frac{1}{1 + e^{-\text{net}_1}} = \frac{1}{1 + e^{-0.3775}} = 0.593$$

$$\text{Out h}_2 = \frac{1}{1 + e^{-\text{net h}_2}} = \frac{1}{1 + e^{-0.393}} = 0.597.$$

$$\text{net O}_1 = \text{out h}_1 * w_5 + \text{out h}_2 * w_6 + b_2$$

$$= (0.593 * 0.4) + (0.597 * 0.45) + 0.6$$

$$= 1.1059$$

$$\text{out O}_1 = \frac{1}{1 + e^{-1.1059}} = 0.751$$

$$= \frac{1}{1 + e^{-1.1059}} = 0.751$$

$$\text{net O}_2 = \text{out h}_1 * w_7 + \text{out h}_2 * w_8 + b_2$$

$$= (0.593 * 0.5) + (0.597 * 0.55) + 0.6$$

$$= 1.2248.$$

~~$$\text{out O}_2 = \frac{1}{1 + e^{-1.2248}} = 0.773.$$~~

~~$$E_1 = \frac{1}{2} (0.01 - 0.751)^2 = 0.275$$~~

~~$$E_2 = \frac{1}{2} (0.99 - 0.773)^2 = 0.02356$$~~

$$\therefore E_{\text{total}} = E_1 + E_2 = 0.275 + 0.02356 = 0.29837$$

$$\begin{aligned} \frac{dE_{\text{total}}}{dw_5} &= \frac{dE_{\text{total}}}{\text{out O}_1} * \frac{\text{out(O)}_1}{\text{out(O)}_1} * \frac{d\text{out(O)}_1}{dw_5} \\ &= -(T_1 - O_1) O_1 (1 - O_1) * \text{out(h)}_1. \\ &= 0.5564 * 0.249 * 0.593 = 0.08215 \end{aligned}$$

$$\begin{aligned} \frac{dE_{\text{total}}}{dw_6} &= -(T_1 - O_1) O_1 (1 - O_1) * \text{out(h)}_2. \\ &= 0.0827 \end{aligned}$$

$$\frac{dG_{\text{total}}}{dW_2} = \frac{dE_{\text{out}}}{dW_2} \times \frac{dW_2}{1 - \alpha_{\text{out}} O_2} \times \frac{d\alpha_{\text{out}}(\text{log})}{dW_2}$$

$$= -(T_2 - O_2) O_2 (1 - O_2) \times \text{out}(h_1).$$

$$= -0.0225 \times 0.773 \times 0.227 \times 0.597$$

$$\frac{dE_{\text{total}}}{dW_2} = -(T_2 - O_2) O_2 (1 - O_2) \times \text{out}(h_2),$$

$$= -0.217 \times 0.773 \times 0.227 \times 0.597$$

$$= -0.0227$$

$$Btu(1) = 1$$

$$Btu(1) = 1$$

$$Btu(1) = 0.125 \times 10^3 + 1 = 13.5 \text{ kJ}$$

$$Btu(1) = 0.125 \times 10^3 + 1 = 13.5 \text{ kJ}$$

$$Btu(1) = 0.125 \times 10^3 + 1 = 13.5 \text{ kJ}$$

13.5 kJ = 13.5 J

13.5 J = 13.5 mJ

13.5 mJ = 13.5 mJ

13.5 mJ = 13.5 mJ

13.5 mJ = 13.5 mJ