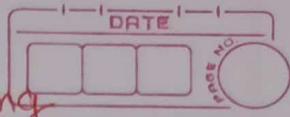


Module 2: Supervised Learning



- * Preference bias: The inductive bias of ID3 is give a preference for certain hypothesis that can be eventually enumerated. This form of bias is called ...
- * Restriction bias: The bias of CANDIDATE-ELIMINATION algo. is in the form of a categorical restriction on the set of hypothesis considered. This form of bias is called ...

Q: Which type of inductive bias is required in order to generalize beyond the training data?

→ Preference bias is more desirable than a restriction bias, because it allows the learner to work within a complete hypothesis space that is assumed to contain the unknown target function.

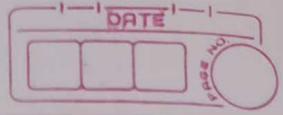
* ISSUES IN DECISION TREE LEARNING:

(a)

→ Avoiding Overfitting the Data

Definition: Given a hypothesis space H , a hypothesis $h \in H$ is said to **overfit** the training data if there exists some alternative -

hypothesis $h' \in H$, such that h has smaller error than h' over the training data, but h' has a smaller error than h over the entire distribution of instances.



⇒ There are several approaches to avoiding overfitting in decision tree learning.

step 2:

(i) Approaches that stop growing the tree earlier, before it reaches the point where it perfectly classifies the training data.

(ii) Approaches that allow the tree to overfit the data, and then post-prune the tree.

Note:

Note:- Although the first of these approaches might seem more direct, the second approach of post-pruning overfit trees has been found to be more successful in practice. This is due to difficulty in the first approach of estimating precisely when to stop growing the tree.

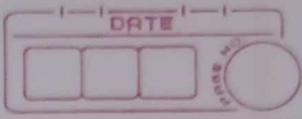
⇒ (i) Reduced Error Pruning:

two different sets like

⇒ We use training set and validation set to prevent overfitting.

step 1:- Reduced Error Pruning is to consider each of the decision nodes in the tree to be candidates for pruning.

Decision



step 2:

Pruning a decision node consists of removing the subtree rooted at that node, making it a leaf node, and assigning it the most common classification of the training examples affiliated with that node.

Note:- Nodes are removed only if the resulting pruned tree performs no worse than the original over the validation set.

- Nodes are pruned iteratively, always choosing the node whose removal most increases the decision tree accuracy over the validation set.
- Pruning of nodes continues until further pruning is harmful (i.e. decreases accuracy of the tree over the validation set).
- In this approach, available data has been split into three subset:
 1. The training examples
 2. The validation examples
 3. The test examples

Drawback:- When the data is limited, withholding part of it for the validation set reduces even further the number of examples available for training.

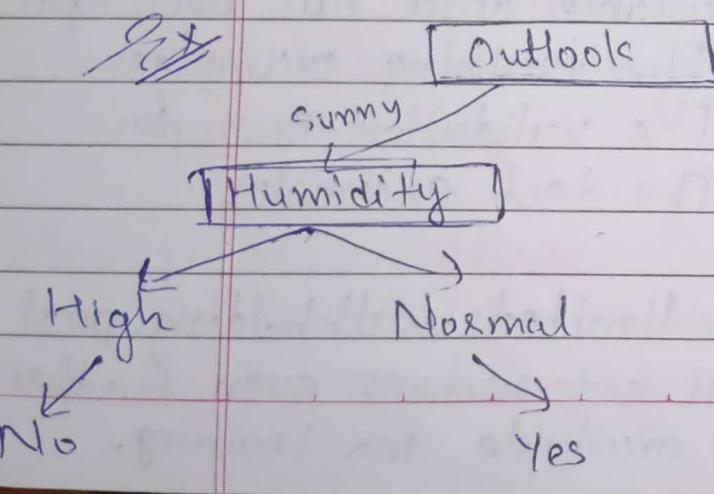
(ii) Rule Post-Pruning:

Step 1: Infer the decision tree from the training set, growing the tree until the training data is fit as well as possible, and allowing overfitting to occur.

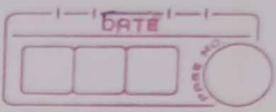
Step 2: Convert the learned tree into an equivalent set of rules by creating one rule for each path from the root node to a leaf node.

Step 3: Prune (generalize) each rule by removing any preconditions that result in improving its estimated accuracy.

Step 4: Sort the pruned rules by their estimated accuracy, and consider them in this sequence when classifying subsequent instances.



~~Rule:-~~
IF (Outlook = Sunny) \wedge (Humidity = High)
THEN PlayTennis = No



★ Why convert the decision tree to rules before pruning?

⇒ Converting to rules allows distinguishing among the different contexts in which a decision node is used. Because each distinct path through the decision tree node produces a distinct rule, the pruning decision regarding that attribute test can be made differently for each path.

⇒ Converting to rules removes the distinction between attribute tests that occur near the root of the tree and those that occur near the leaves.

⇒ Converting to rules improves readability. Rules are often easier for people to understand.

(b) ⇒ Incorporating Continuous-Valued Attributes

⇒ In initial definition of ID3 is restricted to attributes that take on a discrete set of values.

1. Target attribute whose value is predicted by the learned tree must be discrete valued.



2. The attributes tested in the decision nodes of the tree must also be discrete valued.

→ ^{Second} This ~~8~~ restriction can easily be removed so that continuous-valued decision attributes can be incorporated into the learned tree.

→ This can be accomplished by dynamically defining new discrete-valued attributes that partition the continuous attribute value into a discrete set of intervals.

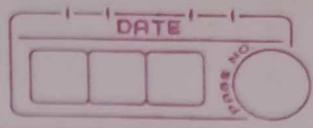
~~Ex~~

Temperature: 40 48 60 72 80 90
PlayTennis : No No Yes Yes Yes No

Step 1:- Pick a threshold, c , that produces the greatest information gain.

→ In the current example, there are two candidate thresholds, corresponding to the values of Temperature at which the value of PlayTennis changes!

$$\frac{(48+60)}{2} = 54 \text{ and } \frac{(80+90)}{2} = 85$$

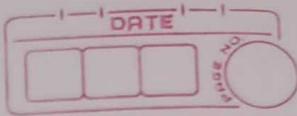


- ⇒ The information gain can then be computed for each of the candidate attributes, Temperature and Temperature_{>85}, and the best can be selected (Temperature_{>85}). ^{>54}
- ⇒ This dynamically created boolean attribute can then compete with the other discrete-valued candidate attributes available for growing the decision tree.

(c) ⇒ Alternative Measures for Selecting Attributes:-

- ⇒ There is a natural bias in the information gain measure that favours attributes with many values over those with few values.
- ⇒ But, simply put, it has so many possible values that it is bound to separate the training examples into very small subsets. Because of this it will have a very high information gain relative to the training examples, despite being a very poor predictor of the target function over unseen instances.

~~Ex~~ Date attribute



Solution: One alternative measure that has been used is gain ratio. The gain ratio measure penalizes attributes such as Date by incorporating a term, called split information, that is sensitive to how broadly and uniformly the attribute splits the data:

$$\text{Split Information}(S, A) = - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

$$\text{Gain Ratio}(S, A) = \frac{\text{Gain}(S, A)}{\text{Split Information}(S, A)}$$

(D) \Rightarrow Handling Training Examples with Missing Attribute Values:-

\Rightarrow One strategy for dealing with the missing attribute value is to assign it the value that is most common among training examples at node n.

\Rightarrow A second procedure is to assign a probability to each of the possible values of A rather than simply assigning the most common value to A(x).

(E) → Handling Attributes with Differing Costs:

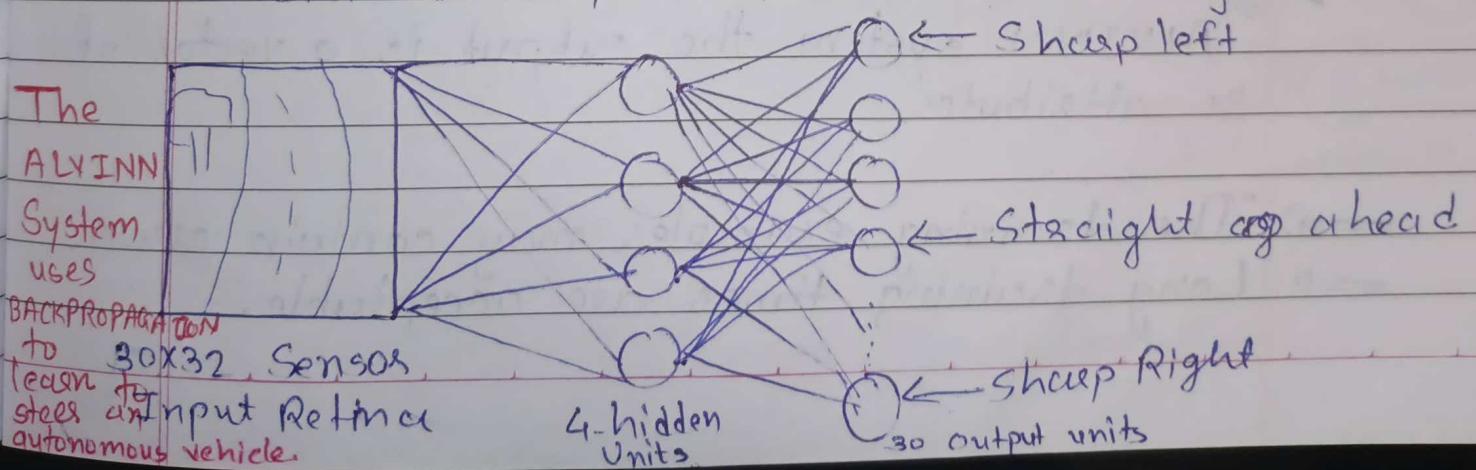
⇒ Measuring attribute cost and prefer cheap cost compare to higher cost attribute. But learner can also use costly/higher cost attribute if it provides good gain value.

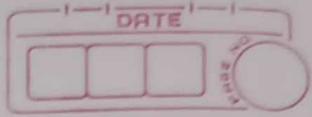
⇒ Problem in mention approach is it is bias towards lower cost attribute.

→ ~~Robot & sonar~~: time required to position medical diagnosis: cost of a laboratory test

Module: 3: Artificial Neural Network & Genetic Algorithms

⇒ Neural network learning methods provide a robust approach to approximating scalar-valued, discrete-valued, and vector-valued target functions.





★ Appropriate Problems for Neural Network Learning:

→ ANN learning is well-suited to problems in which the training data corresponds to noisy, complex sensor data, such as inputs from cameras and microphones.

Note:- The BACKPROPAGATION algorithm is the most commonly used ANN learning technique.

* characteristics of problems where ANN is used:-

→ Instances are represented by many attribute value pairs.

→ The target function output may be discrete-valued, real-valued or a vector of several real or discrete valued attributes.

~~Ex~~

In ALVINN system the output is a vector of 30 attributes.

→ The training examples may contain errors.

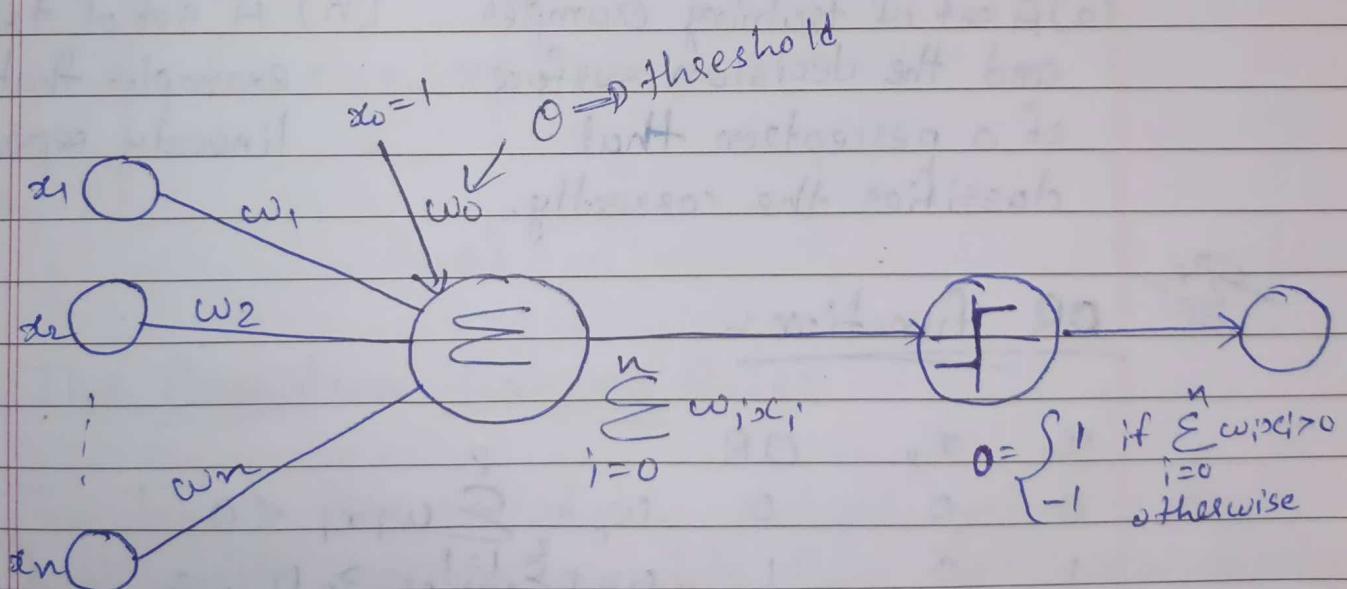
→ Long training times are acceptable.



- Fast evaluation of the learned target function may be required.
- The ability of humans to understand the learned target function is not important.

★ PERCEPTRONS:-

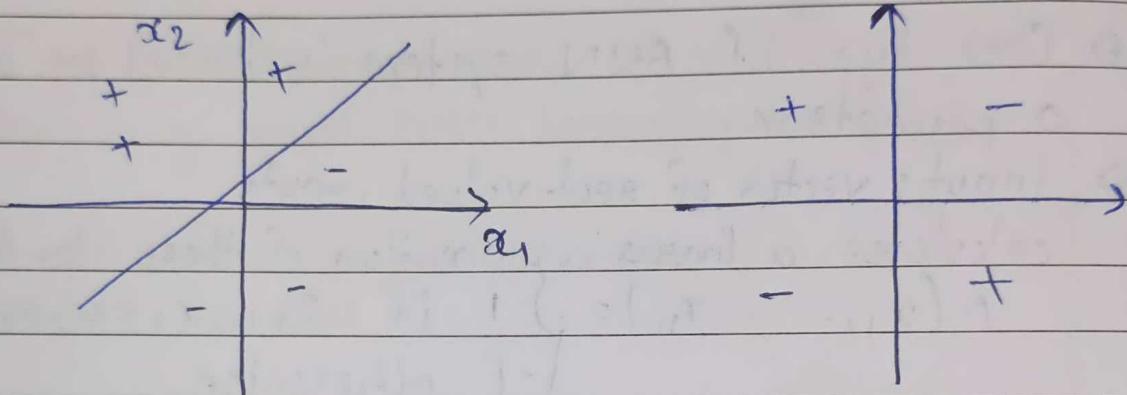
- One type of ANN system is based on a unit called a perceptron.
 - Input: vector of real-valued inputs calculates a linear combination of these i/p's & o/p $\sum_{i=0}^n w_i x_i + b$
- $$O(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$



Representation of Perceptron

→ A single perceptron can be used to represent many boolean functions.

~~Perceptrons~~ Perceptrons can represent all of the primitive boolean functions AND, OR, NAND (\neg AND), and NOR (\neg OR).



(a) A set of training examples and the decision surface of a perceptron that classifies the correctly.

(b) A set of training examples that is not linearly separable.

~~Ex~~ OR function:

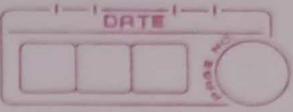
$x_1 \quad x_2 \quad \text{OR}$

$$0 \quad 0 \quad 0 \quad w_0 + \sum_{i=1}^2 w_i x_i < 0$$

$$1 \quad 0 \quad 1 \quad w_0 + \sum_{i=1}^2 w_i x_i > 0$$

$$0 \quad 1 \quad 1 \quad w_0 + \sum_{i=1}^2 w_i x_i > 0$$

$$1 \quad 1 \quad 1 \quad w_0 + \sum_{i=1}^2 w_i x_i \geq 0$$



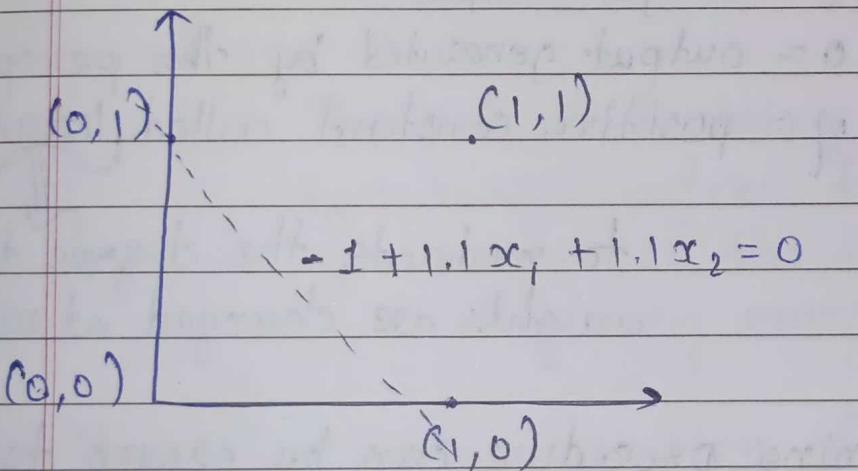
$$w_0 + w_1 \cdot 0 + w_2 \cdot 0 < 0 \Rightarrow w_0 < 0$$

$$w_0 + w_1 \cdot 0 + w_2 \cdot 1 \geq 0 \Rightarrow w_2 \geq -w_0$$

$$w_0 + w_1 \cdot 1 + w_2 \cdot 0 \geq 0 \Rightarrow w_1 \geq -w_0$$

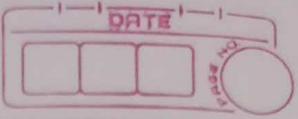
$$w_0 + w_1 \cdot 1 + w_2 \cdot 1 \geq 0 \Rightarrow w_1 + w_2 \geq -w_0$$

\Rightarrow One of the possible solution to this set of inequalities is $w_0 = -1$, $w_1 = 1.1$, $w_2 = 1.1$



* The Perception Training Rule:

\Rightarrow To learn an acceptable weight vector is to begin with random weights, then iteratively apply the perceptron to each training example, modifying the perceptron weights whenever it misclassifies an example.



→ This process is repeated, iterating through the training examples as many times as needed until the perceptron classifies all training examples correctly.

Perceptron training rule;

$$w_i \leftarrow w_i + \Delta w_i$$

where;

$$\Delta w_i = \eta (t - o) x_i$$

t = target output

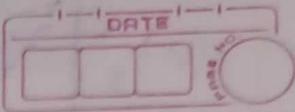
o = output generated by the perceptron

η = positive constant called learning rate

to moderate the degree to which weights are changed at each step.

Note:- Above learning procedure can be proven to converge within a finite number of applications of the perceptron training rule to a weight vector that correctly classifies all training examples, provided the training examples are linearly separable and provided a sufficiently small η is used.

→ If the data are not linearly separable, convergence is not assured.



★ Gradient Descent and the Delta Rule:-

→ Delta rule is designed to overcome the difficulty faced by perceptron rule. That means if the training examples are not linearly separable, the delta rule converges toward a best-fit approximation to the target concept.

$$\text{training error } E(\vec{\omega}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

where;

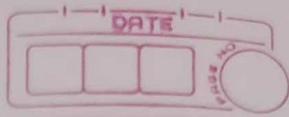
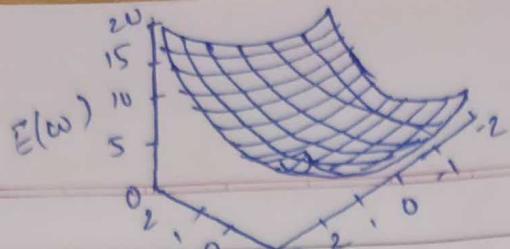
D = set of training examples.

t_d = target output for training example d .

o_d = output of the linear unit for training example d .

Note:- Here, E is a function of $\vec{\omega}$ because the linear unit output o depends on this weight vector.

→ Gradient descent search determines a weight vector that minimizes E by starting with an arbitrary initial weight vector, then repeatedly modifying it in small steps.



→ At each step, the weight vector is altered in the direction that produces the steepest descent (slope) along the error surface. This process continues until the global minimum error is reached.

$$\nabla E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

where,

$\nabla E(\vec{w})$ = vector, whose components are the partial derivatives of E with respect to each of the w_i .

Training rule for gradient descent --

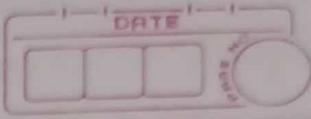
$$\vec{w} = \vec{w} + \Delta \vec{w}$$

where,

$$\Delta \vec{w} = -\eta \nabla E(\vec{w})$$

η = learning rate determines the step size in the gradient descent search.

Note:- negative sign is present because we want to move the weight vector in the direction the decreases E .



* Gradient-Descent algorithm:

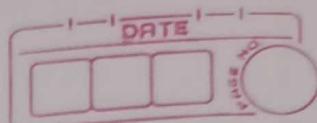
GRADIENT-DESCENT($\text{training_examples}, n$)

Each training example is a pair of the form $\langle \vec{x}, t \rangle$, where \vec{x} is the vector of input values, and t is the target output value.

- Initialize each w_i to some small random value.
- Until the termination condition is met, Do
 - (i) Initialize each Δw_i to zero.
 - (ii) For each $\langle \vec{x}, t \rangle$ in training examples,
 - Input the instance \vec{x} to the unit and compute the output o
 - For each linear unit weight w_i , Do
$$\Delta w_i = \Delta w_i + \eta(t - o)x_i;$$
 - (iii) For each linear unit weight w_i , Do
$$w_i = w_i + \Delta w_i$$

* Practical difficulties in applying gradient descent:

- ⇒ Converging to a local minimum can sometimes be quite slow (i.e. it can require many thousands of gradient descent steps).



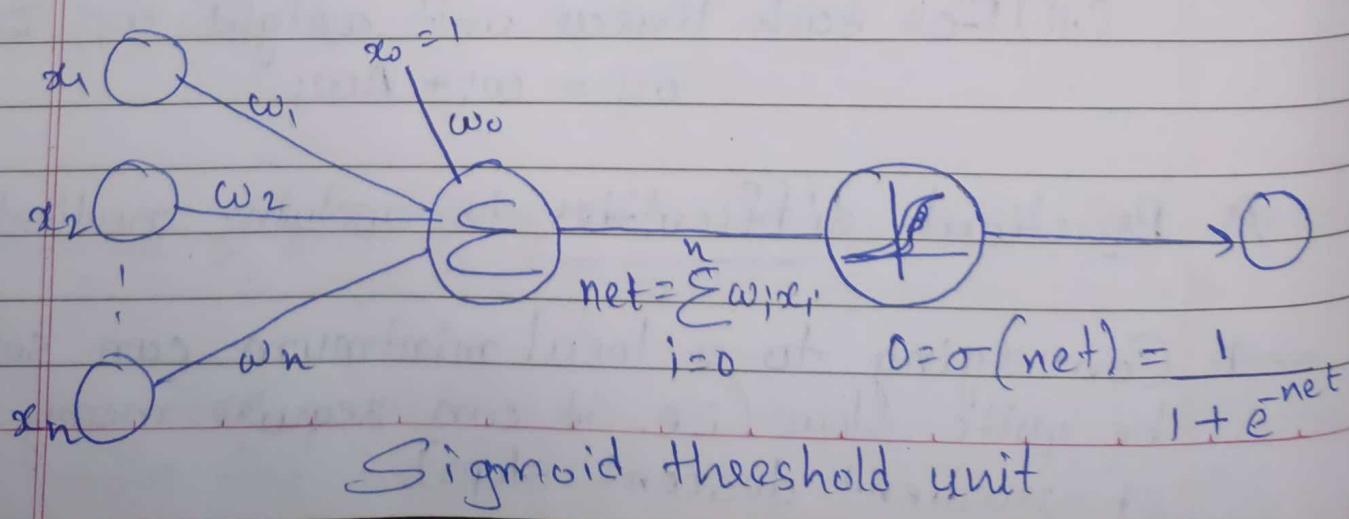
→ If there are multiple local minima in the error surface, then there is no guarantee that the procedure will find the global minimum.

$$\Delta w_i = \eta(t - o) x_i$$

↓
delta rule / least-mean-square rule.

* Multilayer Networks & Backpropagation algo:

- Single perceptrons can only express linear decision surfaces.
- Multilayer networks learned by the Backpropagation algo. are capable of expressing a rich variety of nonlinear decision surfaces.



challenges

- Computational complexity due to a high volume of training data.
- Longer training times
- Human intervention to validate o/p variables.

★ Unsupervised Learning:

- ⇒ Only Inputs
- ⇒ Clustering.
- ⇒ K-means clustering.

★ K-NN algo (Lazy learner algorithm)

- ⇒ Most basic instance-based method.
- ⇒ The nearest neighbors of an instance are defined in terms of the standard euclidean distance. step1 Choose the number of neighbors

Algo: → step2 Calculate Euclidean distance

step3 Identify Nearest Neighbors

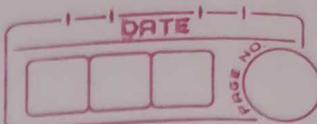
① Training algo:-

- For each training example $\langle x, f(x) \rangle$ add the example to the list training-examples.

② Classification algo:-

- Given a query instance x_q to be classified
- Let x_1, \dots, x_K d -

for weighted K-NN



$\sqrt{d^2}$

Sr. No.	Height	weight	Target	Distance
1	150	50	Medium	8.06
2	155	55	"	2.24
3	160	60	Large	6.71
4	161	59	"	6.40
5	158	65	"	11.05
6	157	54	(?)	

If $k=3$ so 6 will classify in Large

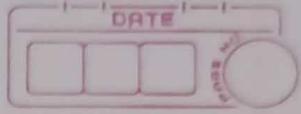
$$d = \sqrt{(157 - 150)^2 + (54 - 50)^2}$$

$$= \sqrt{(7)^2 + (4)^2} = \underline{\underline{8.06}}$$

Supervised

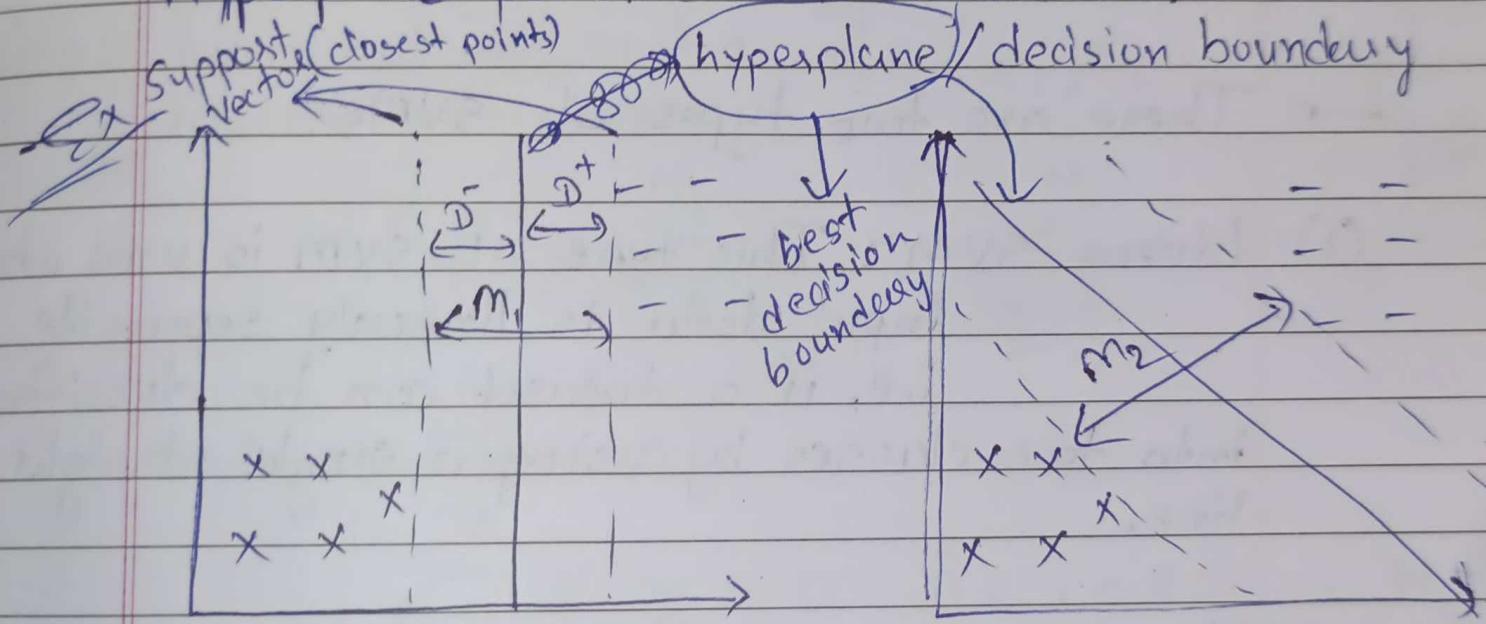
- Linear regression
- SVM
- Naive bayes
- Logistic regression

Unsupervised



★ Support Vector Machine: (Supervised learning)

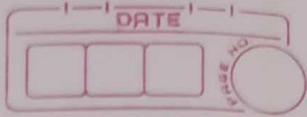
→ To select which decision boundary we have to consider for SVM Maximal Margin Hyperplane parameter is considered.



here; $m_1 \ll m_2$ So second graph. hyper plane is going to be selected for SVM algo

$$\text{Margin} = D^- + D^+$$

Note:- SVM is a classification algorithm and its job is to perform classification by finding a line/plane that classify the data points in a best way.



* SVM Classifies Algorithm

- Supervised learning algorithm
- used to solve both classification & Regression tasks.
- These are two types of SVMs:

(i) Linear SVM: This type of SVM is used when input data is linearly separable, i.e., if a dataset can be classified into two classes by using a single straight line.

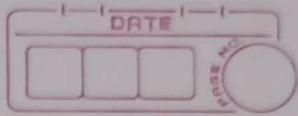
(ii) Non-linear SVM: data is not linearly separable

~~Ex~~

Given a two-class linearly separable dataset of N points of the form :

$$(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_N, y_N)$$

where; y_i 's are either +1 or -1



x_1	x_2	class
2	2	-1
4	5	+1
7	4	+1

Step 1: Find out 'x' vector $\vec{x} = (\alpha_1, \alpha_2, \dots, \alpha_N)$
which maximizes

$$\phi(\vec{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^N \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

No. of training data

dual objective function

pre-conditions:

$$\sum_{i=1}^N \alpha_i y_i = 0 \quad \text{for } i = 1, 2, \dots, N$$

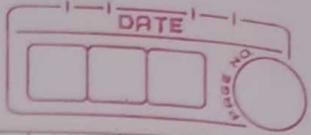
$$\alpha_1 > 0, \alpha_2 > 0, \dots, \alpha_N > 0$$

So

$$-\alpha_1 + \alpha_2 + \alpha_3 = 0$$

Step 2: Find weight vector

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



step 3:- Compute bias

$$b = \frac{1}{2} \left(\min_{i: y_i=+1} (\vec{w} \cdot \vec{x}_i) + \max_{i: y_i=-1} (\vec{w} \cdot \vec{x}_i) \right)$$

↓ only positive example ↓ only negative example

step 4:- The SVM classifier function is given by

$$f(\vec{x}) = \vec{w} \cdot \vec{x} - b$$

↓ eqⁿ of hyperplane

~~Ex~~ Data:-

	x_1	x_2	class
2	1		+1
4	3		-1

here, $N = 2$

$$\vec{x}_1 = (2, 1)$$

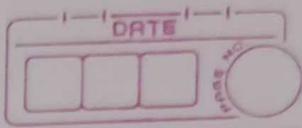
$$\vec{x}_2 = (4, 3)$$

$$y_1 = +1$$

$$y_2 = -1$$

$$f(\vec{x}) = \vec{w} \cdot \vec{x} - b$$

$$\vec{x}_1 \cdot \vec{x}_1 = (2, 1) \cdot (2, 1) \\ = 2 \times 2 + 1 \times 1 \\ = 5$$



Step 1: find $\vec{\alpha} = (\alpha_1, \alpha_2)$

here: $\sum_{i=1}^N \alpha_i y_i = 0$

so,

$$\alpha_1 - \alpha_2 = 0$$

$$\therefore \alpha_1 = \alpha_2$$

and:

$$\alpha_1 > 0 \quad \& \quad \alpha_2 > 0$$

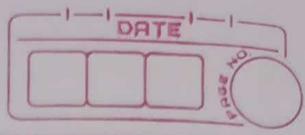
$$\phi(\vec{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^N \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

$$= (\alpha_1 + \alpha_2) - \frac{1}{2} \left[\alpha_1 \alpha_1 y_1 y_1 (\vec{x}_1 \cdot \vec{x}_1) + \alpha_1 \alpha_2 y_1 y_2 (\vec{x}_1 \cdot \vec{x}_2) + \alpha_2 \alpha_1 y_2 y_1 (\vec{x}_2 \cdot \vec{x}_1) \right]$$

$$+ \alpha_2 \alpha_2 y_2 y_2 (\vec{x}_2 \cdot \vec{x}_2) \right]$$

$$= (\alpha_1 + \alpha_2) + \frac{1}{2} \left[\alpha_1^2 (1)(1) (2 \times 2 + 1 \times 1) + \alpha_1 \alpha_2 (1)(-1) \frac{(1)(-1)}{(2 \times 4 + 1 \times 3)} \right.$$

$$+ \alpha_2 \alpha_1 (-1)(+1) (4 \times 2 + 3 \times 1) \left. + \alpha_2^2 (-1)(-1) (4 \times 4 + 3 \times 3) \right]$$



$$= (\alpha_1 + \alpha_2) - \frac{1}{2} [5\alpha_1^2 - 22\alpha_1\alpha_2 + 25\alpha_2^2]$$

→ Find values of α_1 and α_2 which maximizes

$$\phi(\vec{\alpha}) = (\alpha_1 + \alpha_2) - \frac{1}{2} [5\alpha_1^2 - 22\alpha_1\alpha_2 + 25\alpha_2^2]$$

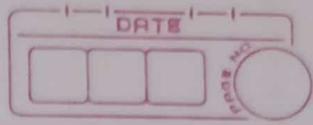
put here; $\alpha_1 = \alpha_2$ So put α , instead of α_2 .

$$\begin{aligned}\phi(\vec{\alpha}) &= (\alpha_1 + \alpha_1) - \frac{1}{2} [5\alpha_1^2 - 22\alpha_1^2 + 25\alpha_2^2] \\ &= 2\alpha_1 - \frac{1}{2} [8\alpha_1^2] \\ &= 2\alpha_1 - 4\alpha_1^2\end{aligned}$$

Note:- If you want to find out maximization of $\phi(\vec{\alpha})$ you have to consider differentiation.

$$\frac{d\phi}{d\alpha_1} = 2 - 8\alpha_1 = 0$$

$$\text{So, } \left[\alpha_1 = \frac{1}{4} = \alpha_2 \right]$$



$$\text{step 2: } \vec{\omega} = \sum_{i=1}^N \alpha_i y_i \vec{x}_i$$

$$= \alpha_1 y_1 \vec{x}_1 + \alpha_2 y_2 \vec{x}_2$$

$$= \frac{1}{4}(+1)(2, 1) + \frac{1}{4}(-1)(4, 3)$$

$$= \frac{1}{4}(2, 1) + \frac{1}{4}(-4, -3)$$

$$= \frac{1}{4} [(2, 1) + (-4, -3)]$$

$$= \frac{1}{4} (-2, -2)$$

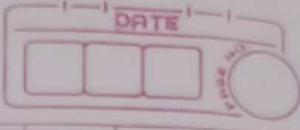
$$= \left(-\frac{1}{2}, -\frac{1}{2} \right)$$

$$\text{step 3: } b = \frac{1}{2} \left(\min_{i: y_i=+1} (\vec{\omega} \cdot \vec{x}_i) + \max_{i: y_i=-1} (\vec{\omega} \cdot \vec{x}_i) \right)$$

$$= \frac{1}{2} ((\vec{\omega} \cdot \vec{x}_1) + (\vec{\omega} \cdot \vec{x}_2))$$

$$= \frac{1}{2} \left[\left(-\frac{1}{2} \times 2 - \frac{1}{2} \times 1 \right) + \left(-\frac{1}{2} \times 4 - \frac{1}{2} \times 3 \right) \right]$$

$$= -\frac{5}{2}$$



step 4: The SVM classifier function is given by

$$f(\vec{x}) = \vec{\omega} \cdot \vec{x} - b$$

$$= \left(-\frac{1}{2}, -\frac{1}{2}\right) \cdot (x_1, x_2) - \left(-\frac{5}{2}\right)$$

$$= -\frac{1}{2}x_1 - \frac{1}{2}x_2 + \frac{5}{2}$$

$$\boxed{= -\frac{1}{2}(x_1 + x_2 - 5)}$$

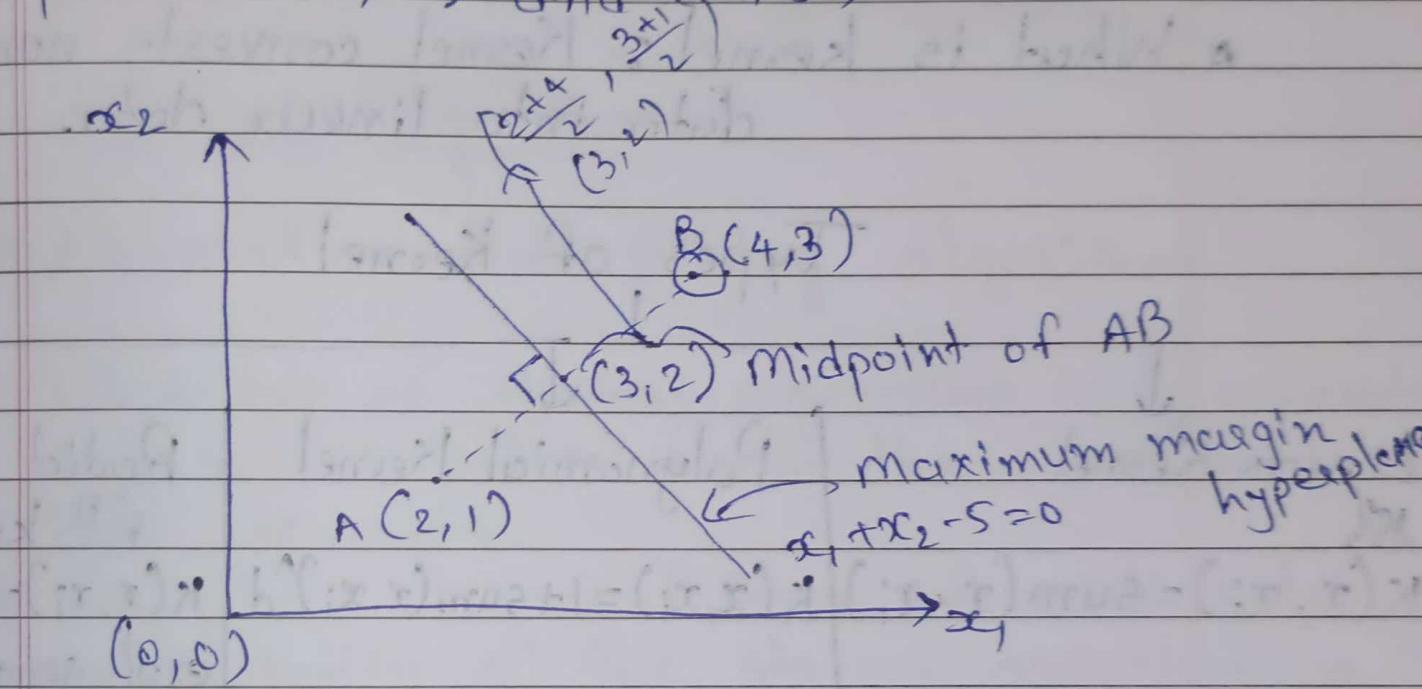
→ The equation of the maximal margin hyperplane is

$$f(\vec{x}) = 0$$

$$\therefore -\frac{1}{2}(x_1 + x_2 - 5) = 0$$

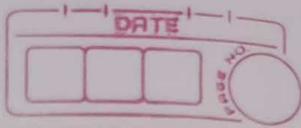
$$\therefore x_1 + x_2 - 5 = 0$$

⇒ The eqⁿ $x_1 + x_2 - 5 = 0$ is the perpendicular bisector of the line segment joining the points $(2, 1)$ and $(4, 3)$



~~Ex/~~ Linear Examples

x_1	x_2	y
3	1	+1
3	-1	+1
6	1	+1
6	-1	+1
1	0	-1
0	1	-1
0	-1	-1
-1	0	-1



Note:-

→ SVM is known for its kernel trick to handle nonlinear data.

* What is kernel :- Kernel converts non-linear data into linear data.

Types of Kernel

Linear Kernel	Polynomial Kernel	Radial Basis Function Kernel
$K(x, x_i) = \sum(x \cdot x_i)$	$K(x, x_i) = 1 + \sum(x \cdot x_i)^d$	$f^n(x, x_i) = \exp(-\gamma \sum((x - x_i)^2))$
(4) Sigmoid Kernel	$f(x_1, x_2) = \tanh(\alpha x_1^T y + b)$	

Imp note:-

high bias \leftarrow low bias :- Before we allowed misclassification we picked a threshold that was very sensitive to the training data.

low variance \rightarrow high variance. It performed poorly when we got new data.



* When we allow misclassifications, the distance between the observations & the threshold is called a Soft margin.

* Apriori Algorithm :-

⇒ Association Rule Mining algorithm.

⇒ Objective of Apriori algorithm is to generate association rules from the frequent itemsets.

⇒ Calculate the confidence of each rule and identify all the strong association rules.

⇒ The association rule describes how two or more objects are related to one another.

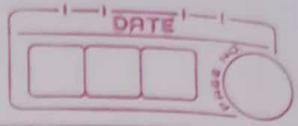
* Support:

- Support refers to the default popularity of any product.

$$\text{Support(item)} = \frac{\text{Trans. of item}}{\text{Total trans.}}$$

* Confidence

- Confidence refers to the possibility that the customers bought multiple items together.



Ex

Transaction ID Items Bought

- | | |
|---|---------------------------------|
| 1 | \$ Bread, Butter, Milk |
| 2 | \$ Bread, Butter |
| 3 | \$ Beer, Cookies, Diapers |
| 4 | \$ Milk, Diapers, Bread, Butter |
| 5 | \$ Beer, Diapers |

Min-support = 40%.

Min-confidence = 70%.

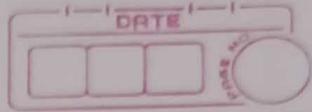
Solⁿ,

Step-1 Unique Item sets are: Bread, Butter, Milk,
Diapers, Beer

1 - Item set	Support count
Bread	3
Butter	3
Milk	2
Beer	2
Cookies	1
Diapers	3

→ show many
time the
item is
bought from
given transaction.

$$\begin{aligned}
 \text{Min-support:count} &= \text{Min-support} \times \frac{\text{no. of trans}}{\text{itemset count}} \\
 &= 40\% \times 5 \\
 &= \frac{40 \times 10^2}{100} \times \frac{8}{28 \times 4} = 2
 \end{aligned}$$



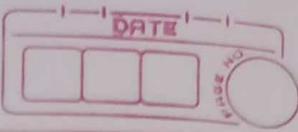
- Min-Support-Count is 2 that indicates that 1-Item set must have support count greater than or equal to 2 then and then it is called 1-frequent Itemset.
- ~~So~~ here cookies has support-count value 1 so we eliminate that itemset

1-Frequent Itemset	Support-count
Bread	3
Butter	3
Milk	2
Beers	2
Diapers	3

step-2: Create combinations for 2-Itemset from ~~1-Frequent Itemsets~~. unique itemset of above table

Ex

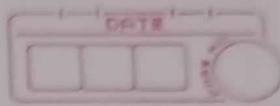
$\{\text{Bread, Butter}\}$, $\{\text{Bread, Milk}\}$, $\{\text{Bread, Beers}\}$,
 $\{\text{Bread, Diapers}\}$, $\{\text{Butter, Milk}\}$, $\{\text{Butter, Beers}\}$,
 $\{\text{Butter, Diapers}\}$, $\{\text{Milk, Beers}\}$,
 $\{\text{Milk, Diapers}\}$, $\{\text{Beers, Diapers}\}$



2- Itemset	Support, count
{Bread, Butter}	3
{Bread, Milky}	2
{Bread, Beer}	0
{Bread, Diapers}	1
{Butter, Milky}	2
{Butter, Beer}	0
{Butter, Diapers}	1
{Milk, Beer}	0
{Milk, Diapers}	1
{Beer, Diapers}	2

2- Frequent Itemset	Support- count
{Bread, Butter}	3
{Bread, Milky}	2
{Butter, Milky}	2
{Beer, Diapers}	2

step:3 Unique items : Bread, Butter, Milk, Beer, Diaper



3-Frequent Itemset	Support-count
Bread, Butter, Milk	2
Bread, Butter, Diapers	1
Bread, Butter, Beer	0
Bread, Milk, Diapers	1
Bread, Milk, Beer	0
Bread, Diapers, Beer	0
Butter, Milk, Diapers	1
Butter, Milk, Beer	0
Butter, Diapers, Beer	0
Milk, Diapers, Beer	0

3-Frequent Itemset	Support Count
Bread, Butter, Milk	2

→ From above 3-frequent Itemset table we can say that it is not possible to generate 4-Itemset because unique itemset available is only 3. So it will stop here and now it will generate association rule.



step:-4 To find strong association rule.

$$\text{min.-confidence} = 70\%$$

$\hookrightarrow \text{Confidence}(x \rightarrow y) = p(y|x) = \frac{p(x \cup y)}{p(x)}$

\hookrightarrow It means this much percent of customers who bought x bought y also.

\Rightarrow here, we have 5 frequent itemsets i.e.

$\{\text{Bread, Butter}\}$, $\{\text{Bread, Milk}\}$, $\{\text{Butter, Milk}\}$
 $\{\text{Diapers, Beer}\}$ and $\{\text{Bread, Butter, Milk}\}$

Therefore, candidate rules are:

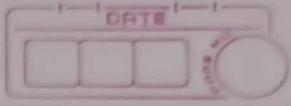
- For $\{\text{Bread, Butter}\}$

$$p(\text{Bread} | \text{Butter}) = \frac{3}{3} = 100\% \text{ (Strong)}$$

$$p(\text{Butter} | \text{Bread}) = \frac{3}{3} = 100\% \text{ (Strong)}$$

- For $\{\text{Bread, Milk}\}$

$$p(\text{Milk} | \text{Bread}) = \frac{2}{3} = 67\%$$



$$P(\text{Bread} \mid \text{Milk}) = \frac{2}{2} = 100\% \text{ (Strong)}$$

- For {Butter, Milk}

$$P(\text{Butter} \mid \text{Milk}) = \frac{2}{2} = 100\% \text{ (Strong)}$$

$$P(\text{Milk} \mid \text{Butter}) = \frac{2}{3} = 67\%$$

\Rightarrow Same way do for {Diaper, Bread} and {Bread, Butter, Milk}.

$$P(\text{Bread, Butter} \mid \text{Milk}) = \frac{2}{2} = 100\% \text{ (Strong)}$$

$$P(\text{Milk, Butter} \mid \text{Bread}) = \frac{2}{3} = 67\%$$

$$P(\text{Milk, Bread} \mid \text{Butter}) = \frac{2}{3} = 67\%$$

$$P(\text{Milk} \mid \text{Bread, Butter}) = \frac{2}{3} = 67\%$$

$$P(\text{Bread} \mid \text{Milk, Butter}) = \frac{2}{2} = 100\% \text{ (Strong)}$$

$$P(\text{Butter} \mid \text{Milk, Bread}) = \frac{2}{2} = 100\% \text{ (Strong)}$$



* Locally Weighted Regression Algorithm

→ Locally weighted regression is instance based learning algorithm.

Local:- Because it will not consider all the data points, it will only consider near data points with the query point.

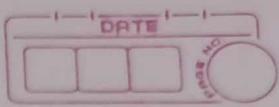
Weighted:- Because the contribution of each training example is weighted by its distance from the query point.

Regression:- Because this is the term used widely in the statistical learning community for the problem of approximating scalar-value functions.

→ Consider, locally weighted regression in which the target function f is approximated near x_q using a linear function of the form ..

query point ↴

$$\hat{f}(x) = w_0 + w_1 q_1(x) + \dots + w_n q_n(x)$$



where;

$a_i(x)$ = value of i^{th} attribute of the instance x .

→ In locally weighted regression initially weights are assigned randomly and to fine tune these weights Gradient Descent algorithm of ANN is used.

$$\text{Error } E(\vec{\omega}) = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

Error f in ANN $d \in D$ target t predicted o

$$E = \frac{1}{2} \sum_{x \in D} [f(x) - \hat{f}(x)]^2$$

Error f in
locally weighted regression

$$\Delta w_j = \eta \sum_{x \in D} [f(x) - \hat{f}(x)] a_j(x)$$

$$w_j = w_j + \Delta w_j$$

weight update rule.

* Locally Weighted Regression algorithms

Approach-1: Locally

Step-1: Minimize the squared errors over just the k-nearest neighbours.

$$E_1(x_q) = \frac{1}{2} \sum_{\substack{x \in \text{k nearest} \\ \text{of } x_q}} (f(x) - \hat{f}(x))^2$$

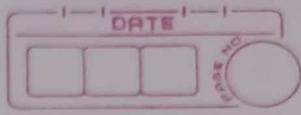
Approach-2: weighted

Minimize the squared errors over the entire set \mathcal{D} of training examples, while weighting the errors of each training example by some decreasing function K of its distance from x_q :

$$E_2(x_q) = \frac{1}{2} \sum_{x \in \mathcal{D}} (f(x) - \hat{f}(x))^2 K(d(x_q, x))$$

Approach-3: Locally weighted (Combine 1 & 2)

$$E_3(x_q) = \frac{1}{2} \sum_{\substack{x \in \text{k nearest} \\ \text{neighbours} \\ \text{of } x_q}} (f(x) - \hat{f}(x))^2 K(d(x_q, x))$$



$$\Delta w_j = \eta \sum_{\substack{x \in k \text{ nearest} \\ \text{neighbors of } x_q}} K(d(x_q, x))(f(x) - \hat{f}(x)) a_j(x)$$

$$w_j = w_j + \Delta w_j$$