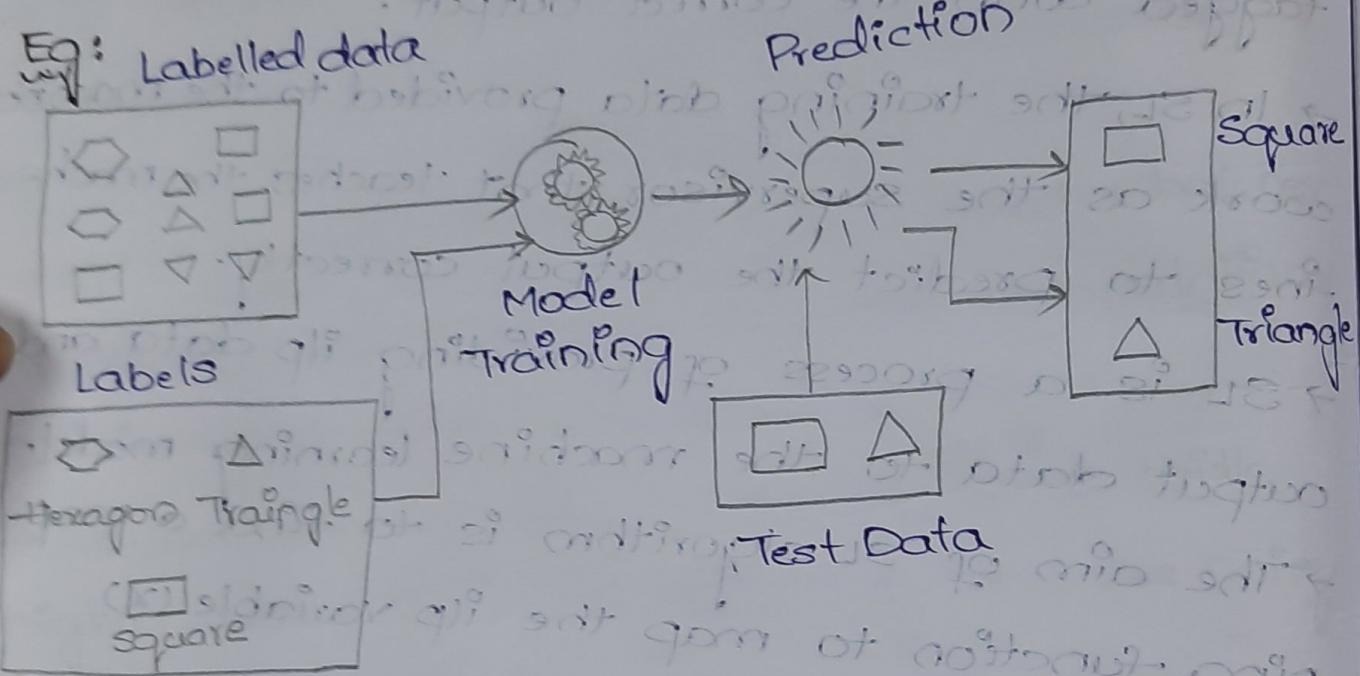


Unit-II

Supervised Learning

- Supervised Learning is one of the types of ML in which machines are trained using 'labelled' training data, and basis of that data, machine predict the output.
- Labelled data means some ILP data is already tagged with the correct output.
- In SL, the training data provided to the machines work as the supervisor that teaches the machines to predict the output correctly.
- SL is a process of providing ILP data and output data to the machine learning model.
- The aim of SL algorithm is to find a mapping function to map the ILP variable (x) with the OLP variable (y).
- In real world SL is used for Risk Assessment, Fraud Detection, Image classification, etc.,
- Spam filtering etc.,

- How supervised learning works?
- In SL models are trained using labelled data set, where model learns about each type of data.
 - Once, training process is completed the model is tested on the basis of test data & then it predicts the output.



- If the given shape has 4 sides, & all are equal, then it will be labelled as a Square.
- If the given shape has 3 sides, then it will be labelled as a Triangle.
- If the given shape has six equal sides then it will be labelled as Hexagon.
- Now, after training, we test our model using test set.

→ The MLG is already trained on all types of shapes, when it finds new shape, it classifies the shape on the basis of no. of sides and predicts the output.

Types of supervised learning algorithms:

Supervised Learning

Regression

Classification

Regression Algorithms: In robotics etc., if there is a relationship between input & output variables.

→ These are used if there is a relationship between input & output variable.

→ It is used for the prediction of continuous variables such as Weather forecasting, Market trends etc.

Some of the popular Regression Alg. are

a) Linear Regression

b) Non-linear Regression

c) Regression Trees

i) Simple Linear Regression ✓

ii) Multiple Regression ✓

iii) Polynomial Regression

iv) Logistic Regression ✓

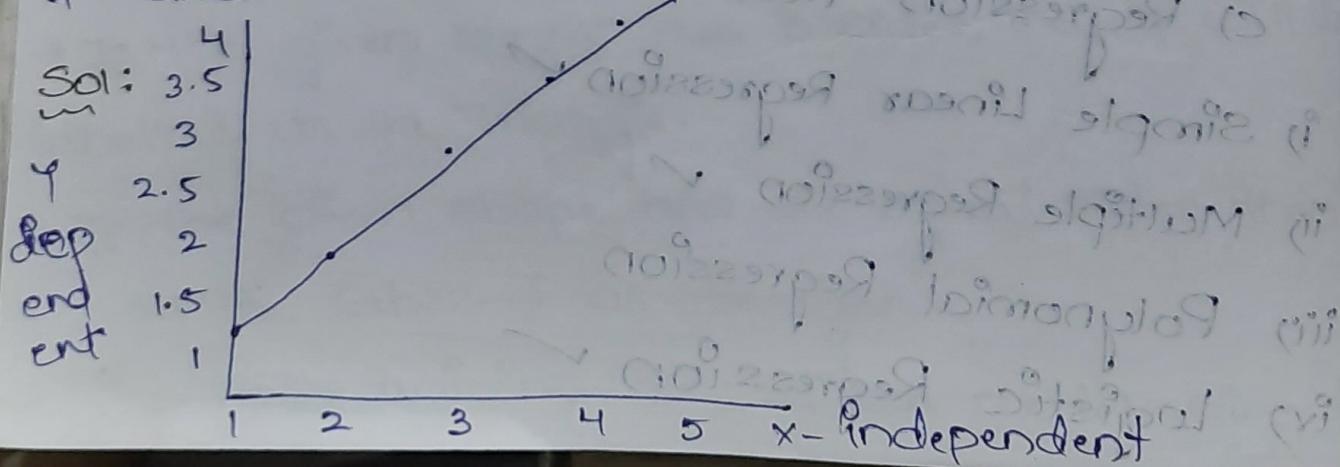
Simple Linear Regression:

→ Assume that there is only one independent variable x , If the relationship b/w x & y (dependent or dep variable) is modeled by the relation, $y = a + bx$, then the regression model is called a Linear Regression Model.

Ex: Let us consider an example, where the five week's sales data (in thousands) is shown.

x_i (Week)	y_i (Sales in thousands)
1	1.2
2	2.6
3	3.2
4	3.8

Apply Linear Regression technique to predict the 7th and 12th week's sales.



$$y = a_0 + a_1 * x + e \quad (e = \text{error})$$

where

$$a_1 = \frac{(\bar{xy}) - (\bar{x})(\bar{y})}{\bar{x}^2 - \bar{x}^2}$$

$a_0 = \text{intercept}$

$a_1 = \text{coefficient of independent variable}$

$\bar{x} = (x)$

$$a_0 = \bar{y} - a_1 * \bar{x}$$

$$\text{sum of } x = 15 \quad \text{Avg} = 3$$

$$\text{sum of } y = 12.6 \quad \text{Avg} = 2.52$$

x_i	y_i	x_i^2	$x_i * y_i$
1	1.2	1	1.2
2	1.8	4	3.6
3	2.6	9	7.8
4	3.2	16	12.8
5	3.8	25	19
sum	15	55	44.4
Average	3	11	8.88

$$\bar{x} = 3, \bar{y} = 2.52, \bar{x}^2 = 11, \bar{xy} = 8.88$$

$$a_1 = \frac{(\bar{xy}) - (\bar{x})(\bar{y})}{\bar{x}^2 - \bar{x}^2} = \frac{(8.88) - (3 * (2.52))}{11 - 3^2}$$

$$= 0.66$$

$$a_0 = \bar{y} - a_1 * \bar{x}$$

$$\bar{x} + a_1 \bar{x} + a_0 = \bar{y}$$

$$\text{Step 2: } \bar{y} = 2.52 - 0.66 * 3 \\ \text{constant} \\ a_0 = 0.54$$

$$\frac{(\bar{y})(\bar{x}) - (\bar{xy})}{\bar{x} - \bar{y}}$$

Linear Regression equation is

$$y = a_0 + a_1 * x$$

$$\bar{x} * a_1 - \bar{y} = a_0$$

$$y = 0.54 + 0.66 * x \quad \text{for } x = 10 \text{ week}$$

The predicted 7th week sale (when $x=7$) is

$$\boxed{\begin{aligned} \hat{y} &= 0.54 + 0.66 * 7 \\ &= 5.16 \end{aligned}}$$

The predicted 12th week sale (when $x=12$) is,

$$\boxed{\begin{aligned} y &= 0.54 + 0.66 * 12 \\ &= 8.46 \end{aligned}}$$

$$22.8 = \bar{y}$$

$$11 = \bar{x}$$

$$62.8 = \bar{P}, 8 = \bar{E}$$

$$22.8 = \bar{y}, 11 = \bar{x}, 62.8 = \bar{P}, 8 = \bar{E}$$

$$\frac{((62.8) * 8) - (22.8)}{62.8 - 11}$$

$$\frac{(\bar{y})(\bar{x}) - (\bar{xy})}{\bar{x} - \bar{y}}$$

$$-10$$

Multiple Linear Regression :

- In linear regression model we have one dependent and one independent variable.
- Multiple regression model involves multiple predictors or independent variables and one dependent variable.
- This is an extension of the linear regression problem.
- The multiple regression of two variables x_1 and x_2 is given as follows :

$$y = f(x_1, x_2)$$

a_0, a_1, a_2 = coefficients of multiple linear regression.

- In general, this is given for 'n' independent variables as

$$y = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_n x_n + \epsilon$$

Here x_1, x_2, \dots, x_n are predictor variables,

y is the dependent variable, $(a_0, a_1, a_2, \dots, a_n)$

are the co-efficients of the regression equation and ϵ is the error term or the residuals of noise.

$$\text{Matrix } \hat{\alpha} = ((x^T x)^{-1} x^T) Y$$

$$x^T x = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{pmatrix} \cdot \begin{pmatrix} 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 2 \end{pmatrix}$$

$4 \times 4 \quad 3 \times 3$

$$= \begin{pmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{pmatrix}$$

3×3

$$\begin{aligned} 1*1+1*1+1*1 &= 3 & 1*1+1*2+1*3+1*4 &= 10 \\ 1*1+2*1+2*1 &= 5 & 1*1+2*2+3*3+4*4 &= 30 \\ 4*1+5*2+8*3+2*4 &= 28 & 4*1+5*2+8*3+2*4 &= 46 \\ 4*1+5*2+8*3+2*4 &= 28 & 4*1+5*2+8*3+2*4 &= 109 \end{aligned}$$

$$(x^T x)^{-1} = \begin{pmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{pmatrix}^{-1}$$

$$= \begin{pmatrix} 3.15 & -0.59 & -0.30 \\ -0.59 & 0.20 & 0.016 \\ -0.30 & 0.016 & 0.054 \end{pmatrix}$$

$$(X^T X)^{-1} X^T = \begin{pmatrix} 3.15 & -0.59 & -0.30 \\ -0.59 & 0.20 & 0.016 \\ -0.30 & 0.016 & 0.054 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{pmatrix} = X^T X$$

$$= \begin{pmatrix} 1.36 & 0.47 & -1.02 & 0.19 \\ 0.05 & 0.47 & -0.098 & 0.155 \\ -0.32 & -0.098 & 0.155 & 0.26 \\ -0.065 & 0.005 & 0.185 & -0.125 \end{pmatrix}$$

$$\rightarrow \hat{a} = ((X^T X)^{-1} X^T) Y$$

$$= \begin{pmatrix} 1.36 & 0.47 & -1.02 & 0.19 \\ 0.05 & 0.47 & -0.098 & 0.155 \\ -0.32 & -0.098 & 0.155 & 0.26 \\ -0.065 & 0.005 & 0.185 & -0.125 \end{pmatrix} \begin{pmatrix} 1 \\ 6 \\ 8 \\ 12 \end{pmatrix}$$

$$= \begin{pmatrix} 1 \\ -1.69 \\ 3.48 \\ -0.05 \end{pmatrix} = a_0, a_1, a_2, a_3 = \hat{a} = ((X^T X)^{-1})$$

\rightarrow Multiple linear regression equation

$$y = a_0 + a_1 x_1 + a_2 x_2$$

$$y = -1.69 + 3.48 x_1 - 0.05 x_2$$

Logistic Regression:

- Linear regression predicts the numerical response but is not suitable for predicting the categorical variables.
- When categorical variables are involved, it is called classification problem.
- Logistic regression is suitable for binary classification problem.
- For example, the following scenarios are instances of predicting categorical variables.
 1. Is the mail spam or not spam?

The answer is Yes or no.
thus, categorical dependent variable is a

binary response of Yes or No.

2. If the student should be admitted or not

is based on entrance exam marks.

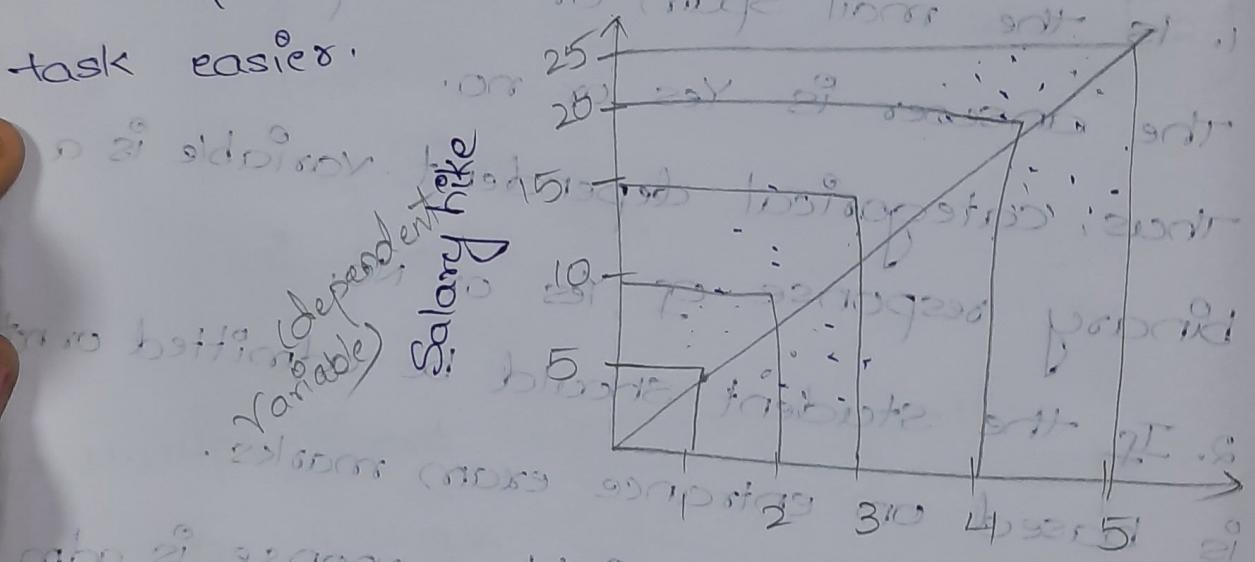
Here, categorical variable response is admitted or not.

3. The student being pass or fail is based on marks secured.

How does the Logistic Regression Algorithm work?

Ex: An organization wants to determine an employee's salary increase based on their performance.

→ For this purpose, a linear regression algorithm will help them decide, considering plotting a regression line by considering the employee's performance as the independent variable, and the salary increase as the dependent variable will make their task easier.



Employee rating

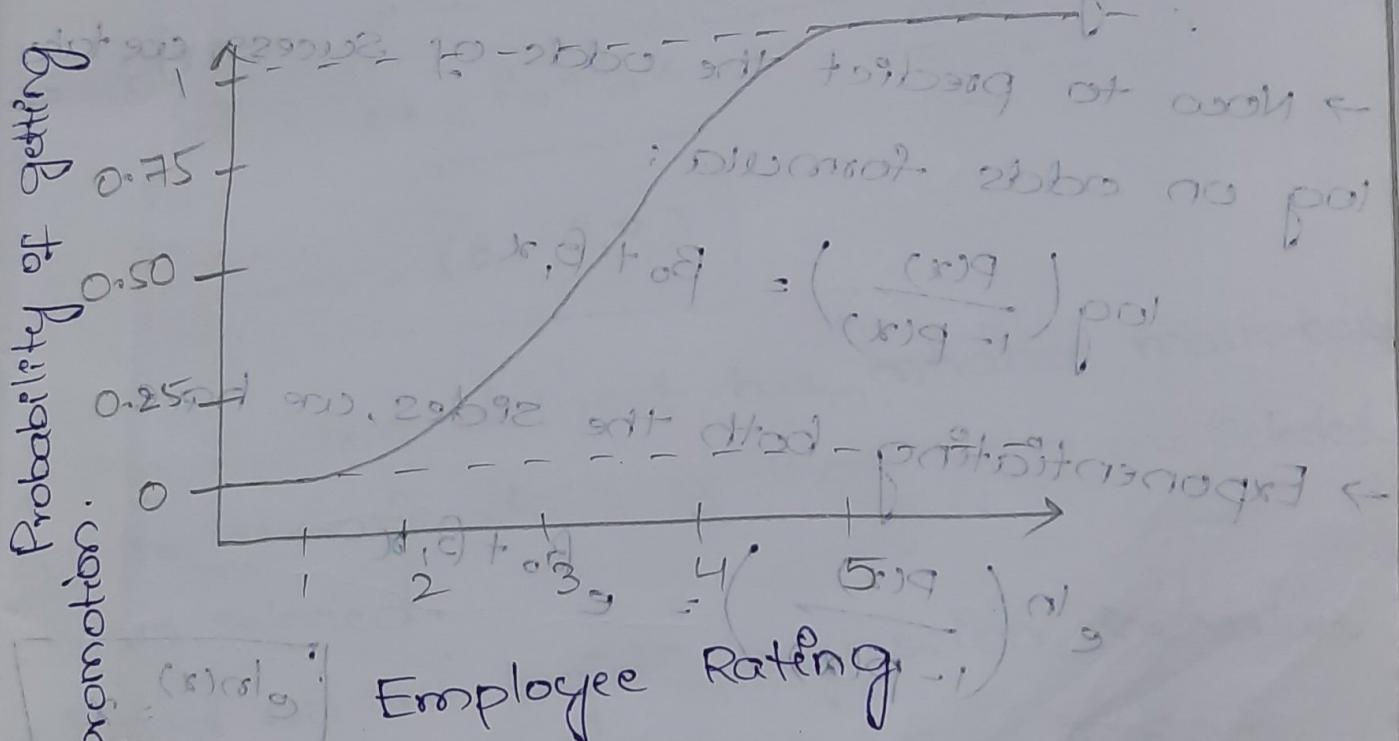
(Independent Variable)

Salary Increase (Dependent Variable)

• Based on the graph, as

- Now, what if the organization wants to know whether an employee would get a promotion or not based on their performance?
- The above linear graph won't be suitable in this case.

- In this case, we clip the line at zero and one.
- As such, we clip the line at zero and one, and convert it into a sigmoid curve (S curve)
- Based on the threshold values, the organization can decide whether an employee will get a salary increase or not.



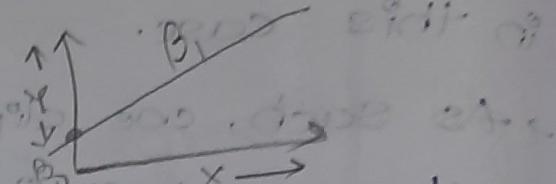
→ the values $\text{bloo}^{(k)}_{\text{og}, i}$ are called 'threshold values'. if $\text{the} \geq 0.5$

To understand logistic regression, let's go over the odds of success.

Probability of an event happening

Odds(Θ) = $\frac{\text{probability of an event happening}}{\text{probability of an event not happening}}$

$$\text{odds}(\Theta) = \frac{P}{1-P}$$



→ odds range from zero to ∞

→ The values of odds range from zero to ∞ and other values of probability lies between 0 & 1.

→ consider the equation of a straight line:

$$y = B_0 + B_1 * x$$

→ Now to predict the odds of success, we take log on odds formula:

$$\log\left(\frac{P(x)}{1-P(x)}\right) = B_0 + B_1 x$$

→ Exponentiating both the sides, we have

$$e^{\ln\left(\frac{P(x)}{1-P(x)}\right)} = e^{B_0 + B_1 x}$$

$$\left(\frac{P(x)}{1-P(x)}\right) = e^{B_0 + B_1 x}$$

$$e^{\ln(x)} = x$$

Formula?

Let $y = e^{\beta_0 + \beta_1 x}$

then $\frac{P(x)}{1 - P(x)} = y$

$$\rightarrow P(x) = y(1 - P(x))$$

$$\rightarrow P(x) = y - y(P(x))$$

$$\rightarrow P(x) + y(P(x)) = y$$

$$\rightarrow P(x)(1+y) = y$$

$$\rightarrow P(x) = \frac{y}{1+y}$$

$$\rightarrow P(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$P(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

Ex: The student dataset has entrance marks based on the historic data of those who are selected or not selected.

→ Based on the logistic regression, the values of the learnt parameters are $\beta_0 = 1$ & $\beta_1 = 80$.

→ Assuming marks of $x = 60$, compute the resultant class.

$$p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} \quad (\beta_0 + \beta_1 x) \text{ for } \\ = \frac{1}{1 + e^{-48}} = 0.44 \quad ((x_1) - 1) \text{ for } x_1 \text{ as } 0.5,$$

→ If we assume the threshold value as 0.5, then it is observed that $0.44 < 0.5$, therefore the candidate with marks $x_1 = 60$ is not selected.

Ex: Hours study

29
15
33
28
39

pass(1) fail(0)

— 0. — (x) 9 ←

β_0

$\beta_1 x_1 + \beta_0$ ←

$\beta_1 x_2 + \beta_0$ ←

$\beta_1 x_3 + \beta_0$ ←

$\beta_1 x_4 + \beta_0$ ←

$\beta_1 x_5 + \beta_0$ ← exam of 5

The dataset of pass or failing an exam of 5

students is given in the table.

use Logistic Regression as classifier to answer the following questions.

1. calculate the probability of pass for students who studied 33 hours.
2. At least how many hours student should study that makes he will pass the course with the probability of more than 95%.

Assume the model suggested by the optimizer for odds of passing the course is

$$\log(\text{odds}) = -64 + 2 * \text{hours}$$

→ We use sigmoid function in logistic regression

$$S(x) = \frac{1}{1+e^{-x}}$$

- 1. calculate the probability of pass for the student who studied 33 hours.

$$S(x) = \frac{1}{1+e^{-x}}$$

(2620.0) pol = (-3) pol

$$z = -64 + 2 * 33$$

$$= -64 + 66 = 2$$

$$\boxed{HP.S = 5}$$

$$\rightarrow P = \frac{1}{1+e^{-2}} = 0.88$$

→ That is, if student studies 33 hours, then

there is 88.1% chance that the student will pass the exam.

$$0.88 = \frac{HP.22}{S} \rightarrow S = 2.44$$

2. At least how many hours student should study that makes he will pass the course with the probability of more than 95% [bbba/pa]

$$P = \frac{1}{1+e^{-z}} = 0.95$$

$$\Rightarrow 0.95 * (1 + e^{-z}) = 1$$

$$\Rightarrow 0.95 * e^{-z} = 1 - 0.95$$

$$\Rightarrow e^{-z} = \frac{0.05}{0.95} \Rightarrow 0.0526,$$

Apply log functions

formula

$$\ln(e^x) = x$$

$$\Rightarrow \log(e^{-z}) = \log(0.0526)$$

$$\Rightarrow -z = \ln(0.0526) = -2.94$$

$$\boxed{z = 2.94}$$

$$\log(\text{odds}) = z = -64 + 2 * \text{hours}$$

$$\Rightarrow 2.94 = -64 + 2 * \text{hours}$$

$$\Rightarrow 2 * \text{hours} = 2.94 + 64$$

$$\Rightarrow 2 * \text{hours} = 66.94$$

$$\Rightarrow \text{hours} = \frac{66.94}{2} = 33.47 \text{ hours.}$$

The student should study at least 33.47 hours so that he will pass the exam with more than 95% probability.

- Linear Regression predicts the numerical response but is not suitable for predicting the categorical variables.
- When categorical variables are involved, it is called classification problem.
- Logistic regression is suitable for binary classification problem.
- Here, the output is often a categorical variable.

Applications of logistic regression:

- Is the mail spam or not spam? The answer is Yes or No. Thus, categorical dependent variable is a binary response of Yes or No.
- If the student should be admitted or not is based on entrance examination marks.
- Here, categorical variable response is admitted (Yes or No).

- Advantages of Linear Logistic Regression:
- Logistic Regression performs better when the data is linearly separable.
 - It does not require too many computational resources as it's highly interpretable.
 - There is no problem scaling the input features.
 - It does not require tuning hyperparameters to implement and train a model using logistic regression.
 - It gives a measure of how relevant a predictor (coefficient size) is, and its direction of association (positive or negative).

Linear Regression

1. Used to solve regression problems.
2. The response variables are continuous in nature.
3. It helps estimate the dependent variable when there is a change in the independent variable.
4. It is a straight line.

Logistic Regression

- Used to solve classification problems.

Particular event taking place based

It is an S-curve
(S = sigmoid).

- Classification - Bayes theorem
- Bayes' Theorem is the cornerstone of Bayesian learning methods because it provides a way to calculate the posterior probability $P(h|D)$ from the prior probability $P(h)$,
- Probability over the data set $P(D)$
 - Current probability $P(D|h)$

$$P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

Maximum A Posteriori (MAP) Hypothesis:
 Maximum A Posteriori (MAP) hypothesis:
 → The learner considers some set of candidate hypotheses H and is interested in finding the most probable hypothesis h given the observed data D or at least one of the maximally probable hypotheses if there are several such maximally probable hypotheses.

→ Any such maximally probable hypothesis is called MAP hypothesis and the MAP hypothesis is determined by

→ We can determine the posterior probability of each candidate hypothesis using Bayes theorem to calculate the posterior probability of each candidate hypothesis.

→ More precisely, one will say that
MAP is a MAP hypothesis provided

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h|D) \text{ where } P(h|D) = \frac{P(D|h)}{P(D)}$$

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \frac{P(D|h) P(h)}{P(D)}$$

Maximum Likelihood and Least-Squared error

hypothesis: ~~MAP~~ (ML) Postulates neural learning approaches such as neural network learning, linear regression & polynomial curve fitting try to learn a continuous-valued target function from training data.

under certain assumptions any learning algorithm that minimizes the squared error loss, the OLP hypothesis provides predictions on the training data.

Maximum Likelihood hypothesis.

The significance of this result is that it provides a Bayesian justification (under certain assumptions) for many neural nets & other curve-fitting methods that attempt to minimize the sum of squared errors over the training data.

→ In order to find the Max likelihood hypothesis
in Bayesian learning for continuous valued
target function, we start the max likelihood hypothesis
definition but using lower case p to refer
to the probability density function.

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} P(D|h)$$

→ We assume a fixed set of training instances

(x_1, \dots, x_m) and therefore consider the data D to
be the corresponding sequence of target values

$$D = (d_1, \dots, d_m) \quad \text{samples =}$$

→ Here, we can write $P(D|h)$ as the product
of the various probabilities:

$$\underset{i=1}{\operatorname{argmax}} P(d_i|h)$$

softmax function or softmax function is used
in statistics and machine learning. It is a
function that takes a vector of real numbers
and normalizes it into a probability distribution.

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

softmax function is also known as a soft-

Naïve Bayes classifier: Algorithm:
 → Bayes theorem is the cornerstone of Naïve Bayes classifier because it provides a way to calculate the posterior probability $P(h|D)$, from the prior probability $P(h)$, together with $P(D)$ and $P(D|h)$.

$$P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

$$\text{MAP} = \underset{h \in H}{\operatorname{argmax}} P(h|D)$$

$$= \underset{h \in H}{\operatorname{argmax}} \frac{P(D|h) P(h)}{P(D)}$$

$$= \underset{v_j \in V}{\operatorname{argmax}} P(D|v_j) P(v_j)$$

→ The Bayesian approach to classifying the new instance is to assign the most probable target value, v_{MAP} given the attribute values $\{a_1, a_2, \dots, a_n\}$ that describe the instance.

$$v_{MAP} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j | a_1, a_2, \dots, a_n)$$

→ We can use Bayes theorem to rewrite this expression as

$$V_{MAP} = \underset{v_j \in V}{\operatorname{argmax}} \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

$$= \underset{v_j \in V}{\operatorname{argmax}} P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

Naive Bayes classifier:

$$V_{NIB} = \underset{v_j \in V}{\operatorname{argmax}} v_j P(v_j) \prod_i P(a_i | v_j)$$

	Day	Outlook	Temp	Humidity	Wind	Play Tennis
D ₁	Sunny	Hot	High	Weak	No	
D ₂	Sunny	Hot	High	Strong	No	
D ₃	Overcast	Hot	High	Weak	Yes	
D ₄	Rain	Mild	High	Weak	Yes	
D ₅	Rain	Cool	Normal	Weak	Yes	
D ₆	Rain	Cool	Normal	Strong	No	
D ₇	Overcast	Cool	Normal	Strong	Yes	
D ₈	Sunny	Mild	High	Weak	No	
D ₉	Sunny	Cool	Normal	Weak	Yes	
D ₁₀	Rain	Mild	Normal	Weak	Yes	
D ₁₁	Sunny	Mild	Normal	Strong	Yes	
D ₁₂	Overcast	Mild	High	Strong	Yes	
D ₁₃	Overcast	Hot	Normal	Weak	Yes	
D ₁₄	Rain	Mild	High	Strong	No	

→ New instance to be classified

outlook = sunny, temp = cool, humidity

(iv) Humidity = high, Wind = strong >

→ Naive Bayes classifier

$$y_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j | T) P(a_j | v_j)$$

$$P(\text{yes} | \text{new instance}) = P(\text{yes}) * P(\text{sunny} | \text{yes}) *$$

$$P(\text{no} | \text{new instance}) = P(\text{no}) * P(\text{sunny} | \text{no}) *$$

$$P(\text{yes}) = \frac{10}{14} \quad \text{and} \quad P(\text{no}) = \frac{4}{14}$$

$$P(\text{yes} | \text{new instance}) = P(\text{yes}) * P(\text{sunny} | \text{yes}) *$$

$$P(\text{no} | \text{new instance}) = P(\text{no}) * P(\text{sunny} | \text{no}) *$$

$$P(\text{yes}) = \frac{10}{14} \quad \text{and} \quad P(\text{no}) = \frac{4}{14}$$

$$P(\text{yes}) = \frac{10}{14} \quad \text{and} \quad P(\text{no}) = \frac{4}{14}$$

	Outlook	Yes	No
Normal	Sunny	219	315
Windy	Overcast	419	105
Rainy		319	315

	Humidity	Yes	No
High	Normal	319	415
Normal	High	619	115

Temp	Yes	No
Hot	219	315
Mild	419	215
Cool	319	215

Wind	Yes	No
Strong	319	315
Weak	619	215

$$P(\text{yes}|\text{new instance}) = 0.0053$$

$$P(\text{no}|\text{new instance}) = 0.0206$$

As $P(\text{no}|\text{new instance})$ is more, hence, the new example is classified as 'No'.

$$VNB(\text{Yes}) = \frac{VNB(\text{Yes})}{VNB(\text{Yes}) + VNB(\text{No})} = 0.205$$

$$VNB(\text{No}) = \frac{VNB(\text{No})}{VNB(\text{Yes}) + VNB(\text{No})} = 0.795$$

Naive Bayes classifier - Example 2:

→ Estimate conditional probabilities of each attribute { color, legs, height, smelly } for

the species classes:

$P(M) = \frac{4}{8} = 0.5$ of the

$P(H) = \frac{4}{8} = 0.5$ of the

table.

→ Using these probabilities estimate the probability values for the new instance -

(color = Green, legs = 2, height = tall and

smelly = No).

$P(M) = 0.5$

$P(H) = 0.5$

$P(L=2) = 0.5$

$P(H=tall) = 0.5$

$P(S=\text{No}) = 0.5$

No	color	Legs	Height	Smelly	Species
1	White	3	short	Yes	M
2	Green	2	Tall	NO	M
3	Green	3	short	Yes	M
4	White	3	short	Yes	M
5	Green	2	short	NO	H
6	White	2	tall	NO	H
7	White	2	Tall	NO	H
8	White	2	short	Yes	H

New Instance = (color=Green, legs=2, height=Tall follows up and smelly= NO)

$$P(M) = \frac{4}{8} = 0.5$$

$$P(H) = \frac{4}{8} = 0.5$$

color	M	H
White	2/4	3/4
Green	2/4	1/4

height	M	H
Tall	3/4	2/4
Short	1/4	2/4

Legs	M	H
2	1/4	4/4
3	3/4	0/4

Smelly	M	H
Yes	3/4	1/4
NO	1/4	3/4

$$P(M | \text{New Instance}) = P(M) * P(\text{color} = \text{Green} | M) \\ * P(\text{Legs} = 2 | M) * P(\text{Height} = \text{tall} | M) * P(\text{smelly} \\ = \text{no} | M)$$

$$\Rightarrow P(M | \text{New Instance}) = 0.5 * \frac{3}{4} * \frac{1}{4} * \frac{3}{4} * \frac{1}{4} \\ = 0.017$$

$$P(H | \text{New Instance}) = P(H) * P(\text{color} = \text{Green} | H) \\ * P(\text{Legs} = 2 | H) * P(\text{Height} = \text{tall} | H) * P(\text{smelly} = \text{no} | H)$$

$$\Rightarrow P(H | \text{New Instance}) = 0.5 * \frac{1}{4} * \frac{4}{4} * \frac{2}{4} * \frac{3}{4} \\ = 0.047$$

$$P(H | \text{New Instance}) > P(M | \text{New Instance})$$

Instance belongs to species

Hence,

Example 3:

s.no	color	Type	Origin	stolen
1	Red	sports	Domestic	Yes
2	Red	sports	Domestic	No
3	Red	sports	Domestic	Yes
4	Yellow	sports	Domestic	No
5	Yellow	sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	sports	Imported	Yes.

New instance = (Red, SUV, Domestic)

$$P(\text{Yes}) = \frac{5}{10} = 0.5 ; P(\text{No}) = \frac{5}{10} = 0.5$$

Color	Yes	No
Red	2/5	2/5
Yellow	2/5	3/5

Type	Yes	No
sports	4/5	2/5
SUV	1/5	3/5

Origin	Yes	No
Domestic	2/5	3/5
Imported	3/5	2/5

$$P(Y|n) = P(\text{Yes}) * P(\text{color=Red|Yes}) * P(\text{Type=SUV|Yes}) * P(\text{or=Domestic})$$

$$= \frac{5}{10} * \frac{3}{5} * \frac{1}{5} * \frac{2}{5} = \frac{3}{125} = 0.024$$

$$P(N|n) = P(\text{No}) * P(\text{color=Red|No}) * P(\text{Type=SUV|No}) * P(\text{or=Domestic})$$

$$= \frac{5}{10} * \frac{2}{5} * \frac{3}{5} * \frac{3}{5} = \frac{9}{125} = 0.072$$

NO

Example 3:

S.no	color	Type	Origin	stolen
1	Red	sports	Domestic	Yes
2	Red	sports	Domestic	No
3	Red	sports	Domestic	Yes
4	Yellow	sports	Domestic	No
5	Yellow	sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	sports	Imported	Yes

New instance = (Red, SUV, Domestic)

$$P(\text{Yes}) = \frac{5}{10} = 0.5 \quad ; \quad P(\text{No}) = \frac{5}{10} = 0.5$$

color	Yes	No
Red	2/5	2/5
Yellow	2/5	3/5

Type	Yes	No
sports	4/5	2/5
SUV	1/5	3/5

Origin	Yes	No
Domestic	2/5	3/5
Imported	3/5	2/5

$$P(Y|n) = P(\text{Yes}) * P(\text{color=Red|Yes}) * P(\text{Type=SUV|Yes}) * P(\text{or=Domestic})$$

$$= \frac{5}{10} * \frac{3}{5} * \frac{1}{5} * \frac{2}{5} = \frac{3}{125} = 0.024$$

$$P(N|n) = P(\text{No}) * P(\text{color=Red|No}) * P(\text{Type=SUV|No}) * P(\text{or=Domestic})$$

$$= \frac{5}{10} * \frac{2}{5} * \frac{3}{5} * \frac{3}{5} = \frac{9}{125} = 0.072$$

NO

K-Nearest Neighbour Classifier:

- KNN is based on feature similarity, we can do classification (is done using distance measure of new instance) using KNN classifier.
- Def: KNN is one of the simplest supervised ML algorithm mostly used for classification. It is also called as Lazy learning.

Example 1: Given data Query = x constant
 $x = (\text{maths} = 6, \text{CS} = 8)$ & $K=3$ nearest neighbour
classification = pass / fail

Maths	CS	Result
4	3	F
6	8	P
7	8	P
5	5	F
8	8	P

Euclidean Distance (d)

$$d = \sqrt{(x_{01} - x_{A1})^2 + (x_{02} - x_{A2})^2}$$

$o \rightarrow$ observed value

$a \rightarrow$ actual value.

$$1) \text{ calculate } d_1 = \sqrt{(6-4)^2 + (8-3)^2} \\ = \sqrt{2^2 + 5^2} = \sqrt{29} = 5.38$$

$K=3$

$$2) \text{ calculate, } d_2 = \sqrt{(6-6)^2 + (8-7)^2} = \sqrt{1+1} = 1$$

$$3) d_3 = \sqrt{(6-7)^2 + (8-8)^2} = \sqrt{1+0} = 1$$

$$4) d_4 = \sqrt{(6-5)^2 + (8-5)^2} = \sqrt{1+9} = \sqrt{10} = 3.16$$

$$5) d_5 = \sqrt{(6-8)^2 + (8-8)^2} = \sqrt{4+0} = \sqrt{4} = 2$$

compare with 3 neighbours.

$d_2, d_3, d_5 = \text{Pass}$

Instance-based learning Ex: x_1, x_2, x_3, x_4, x_5
 value can be predicted based on the previous examples by calculating the distance measure.

Ex: S.No	Height	Weight	Target
1	150	50	Medium
2	155	55	Medium
3	160	60	Large
4	161	59	Large
5	158	65	Large
	157	54	?

Given $K=3$, $d_1=8.06$, $d_2=2.24$, $d_3=6.71$

$d_4=6.40$, $d_5=11.05$ $[d_6 = \text{Large}]$

→ Ex 2 KNN

Height (cm)	Weight (kg)	Class
167	51	underweight
182	62	Normal
176	69	Normal
173	64	Normal
172	65	Normal
174	56	underweight
169	58	Normal
173	57	Normal
170	55	Normal
170 (x_2)	57 (y_2)	?

(x_2, y_2 - fixed)

First find the nearest neighbours.

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

$$d_1 = \sqrt{(170 - 167)^2 + (57 - 51)^2}$$

$$d_2 = \sqrt{(170 - 182)^2 + (57 - 62)^2}$$

$$d_1 = 6.7; d_2 = 13; d_3 = 13.4; d_5 = 7.6;$$

$$d_6 = 8.2; d_7 = 4.1; d_8 = 1.4; d_9 = 3$$

→ Sort based on distance.

If $K=1$ then it is Normal.

If $K=2$, then Normal.

KNN-3:

The Restaurant A sells burger with optional flavor pepper, Ginger, and chilly.
Everyday this week you have tried a burger at A
and kept a record of which you liked.
using Hamming distance, show how the KNN classifier
with majority voting could classify.

{ Pepper: false, ginger = true, chilly = true }

Day	Pepper	Ginger	Chilly	Liked	Distance
②-A	True	True	True	False	$H=0+0=1$
B	True	False	False	True	$H=1+1=3$
①-C	False	True	True	False	$H=0+0=0$
③-D	False	True	False	True	$H=0+0+1=1$
E	True	False	False	True	$H=1+1=3$

→ In Hamming distance, if the values are same
that is x_1 is same as x_2 then distance is 0 or,
ex: If x_1 = blue; x_2 = blue then distance b/w x_1 & x_2 = 0

If x_1 = blue; x_2 = red then $x_1 \neq x_2 = 1$

→ From this 3, we have to find new table.

Majority is False.
so, new example is classified as False in

this case.

Support Vector Machine: Introduction

- SVM is one of the most popular supervised learning algorithms, which is used for classification as well as regression problems.
- primarily, it is used for classification problems in ML.
- The goal of the SVM algorithm is to create the best line (or) decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future.
- The best decision boundary is called Hyperplane.
- How to draw Hyperplane?
- The dimensions of the hyperplane depend on the features present in the dataset, which means if there are 2 features, then hyperplane will be a straight line.
- If there are 3 features, then hyperplane will be a 2-dimensional plane.
- We always create a hyperplane that has a maximum margin, which means the maximum distance between the data points.

SVM can be two types:

1. Linear SVM

2. Non-linear SVM

1. Linear SVM: Linear SVM is used for linearly separable data, which means if a dataset can be classified in to two classes by using a single straight line, then such data is termed as linearly separable data, and classifier is used called as linear SVM classifier.

2. Non-linear SVM: It is used for non-linearly separated data, which means if a dataset cannot be classified by using a straight line then such data is termed as non-linear data and classifier used is called as Non-linear SVM classifier.

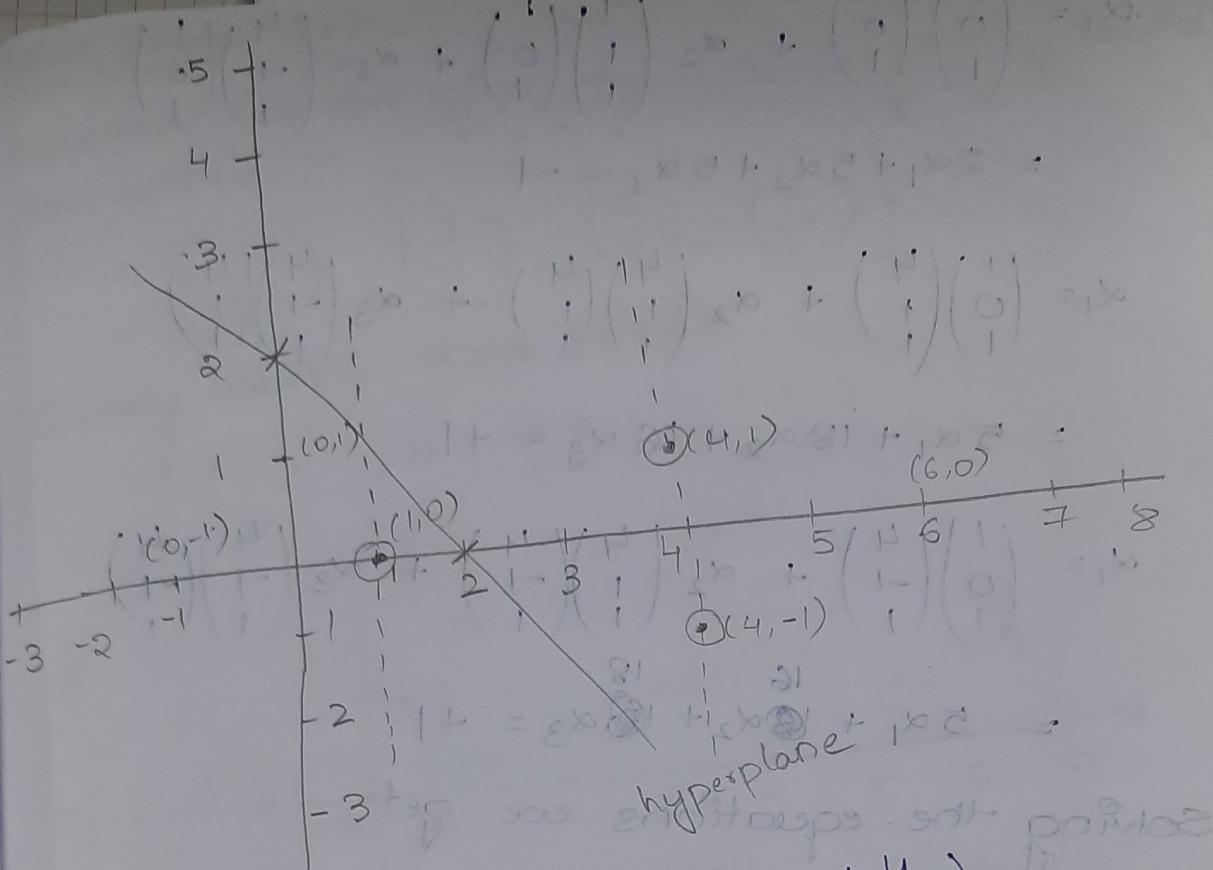
SVM - Example:

→ points $(4, 1)$, $(4, -1)$ and $(6, 0)$ belong to class +ve

→ points $(1, 0)$, $(0, 1)$ & $(0, -1)$ belong to -ve.

→ Draw an optimal hyperplane to classify the points.

→ Margin is the distance between minimum & maximum distance from the classifier.



$s_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}; s_2 = \begin{pmatrix} 4 \\ 1 \end{pmatrix}; s_3 = \begin{pmatrix} 4 \\ -1 \end{pmatrix}$

The augmented vector can be obtained by adding the bias given as follows:

$$\bar{s}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}; \bar{s}_2 = \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix}; \bar{s}_3 = \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix}$$

set of 3 equations can be obtained based on 3 SVM

$$\alpha_1 \bar{s}_1 \bar{s}_1 + \alpha_2 \bar{s}_2 \bar{s}_1 + \alpha_3 \bar{s}_3 \bar{s}_1 = 1$$

$$\alpha_1 \bar{s}_1 \bar{s}_2 + \alpha_2 \bar{s}_2 \bar{s}_2 + \alpha_3 \bar{s}_3 \bar{s}_2 = 1$$

$$\alpha_1 \bar{s}_1 \bar{s}_3 + \alpha_2 \bar{s}_2 \bar{s}_3 + \alpha_3 \bar{s}_3 \bar{s}_3 = 1$$

$$\alpha_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

$$= 2\alpha_1 + 5\alpha_2 + 5\alpha_3 = -1$$

$$\alpha_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix}$$

$$= 5\alpha_1 + 18\alpha_2 + 16\alpha_3 = +1$$

$$\alpha_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix}$$

$$= 5\alpha_1 + 16\alpha_2 + 18\alpha_3 = +1$$

Solving the equations we get

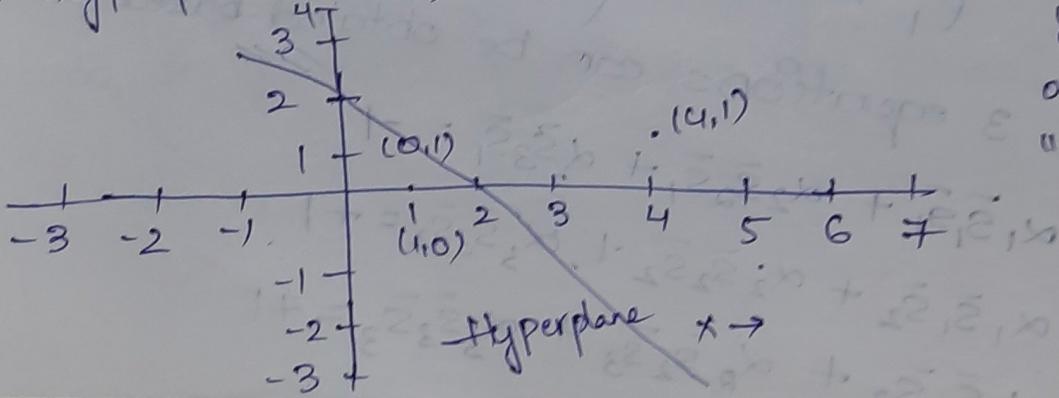
$$\alpha_1 = -3; \alpha_2 = +1; \alpha_3 = 0.$$

The optimal hyperplane is given as

$$w = \sum_{i=1}^3 \alpha_i * \bar{s}_i$$

$$= -3 * \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 1 * \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix} + 0 * \begin{pmatrix} 4 \\ -1 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ -2 \end{pmatrix}$$

The hyperplane is $(1, 1)$ with an offset -2.



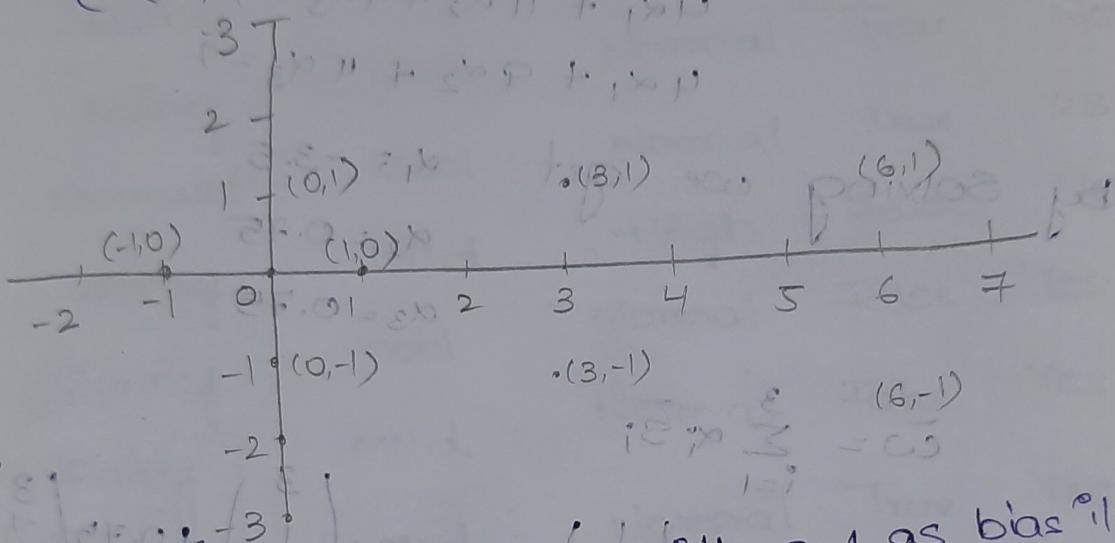
SVM - Linear - Ex:

Suppose we are given the following +ve/-ve labeled data points,

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

a -ve/-ve labeled data points

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



Each vector is augmented (with a 1 as bias).
so, $s_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$ then $\bar{s}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$

$s_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix}$ then $\bar{s}_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix}$ & $s_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$ then $\bar{s}_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$

3 equations:

$$\alpha_1 \bar{s}_1 \cdot \bar{s}_1 + \alpha_2 \bar{s}_2 \cdot \bar{s}_1 + \alpha_3 \bar{s}_3 \cdot \bar{s}_1 = +1$$

$$\alpha_1 \bar{s}_1 \cdot \bar{s}_2 + \alpha_2 \bar{s}_2 \cdot \bar{s}_2 + \alpha_3 \bar{s}_3 \cdot \bar{s}_2 = +1$$

$$\alpha_1 \bar{s}_1 \cdot \bar{s}_3 + \alpha_2 \bar{s}_2 \cdot \bar{s}_3 + \alpha_3 \bar{s}_3 \cdot \bar{s}_3 = +1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} = 1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = 1$$

by solving $\left. \begin{array}{l} 2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1 \\ 4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1 \\ 4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1 \end{array} \right\}$

by solving we get $\alpha_1 = -3.5$
 $\alpha_2 = 0.75$
 $\alpha_3 = 0.75$

$$\bar{\omega} = \sum_{i=1}^3 \alpha_i \vec{s}_i$$

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

Finally, remembering that our vectors are aligned with a bias.

→ We equate the last entry in $\bar{\omega}$ as the hyperplane offset b and write the separating,

hyperplane equation, $y = \omega \cdot x + b$

$$\omega = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, b = -2$$

$$1 + \sqrt{2}x_1 + \sqrt{2}x_2 + \sqrt{2}x_3 + \sqrt{2} \text{ hyperplane}$$

CART Classification				Regression Techniques		
Ex:	outlook	Temp	Humidity	Windy	play	
	sunny	Hot	High	False	No	
	sunny	Hot	High	True	No	
	Overcast	Hot	High	False	Yes	
	Rainy	Mild	High	False	Yes	
	Rainy	Cool	Normal	False	No	
	Rainy	Cool	Normal	True		
	Rainy	Cool	Normal	True	Yes	
	Overcast	Mild	High	False		
	Sunny	Mild	Normal	False	Yes	
	Sunny	Cool	Normal	False	Yes	
	Rainy	Mild	Normal	True	Yes	
	Sunny	Mild	High	True	Yes	
	Overcast	Mild	Normal	False	Yes	
	Overcast	Hot	Normal	False	No	
	Rainy	Mild	High	True	No	

→ For the given play Tennis Data set apply the Decision Tree algorithm & find the optimal decision tree.

→ Also predict class label for the following example

Sunny	Hot	Normal	True	?
-------	-----	--------	------	---

Outlook	Overcast	4	Yes	4
	Sunny	5	Yes	2
	Rainy	5	No	3

Attribute	Rules	Error	Total Error
outlook	sunny → No	2/5	4/14
	overcast → Yes	0/4	
	Rainy → Yes	2/5	

Temp	Rules	Error	Total Error
Temp	Hot → No	2/4	5/14
	Mild → Yes	4/6	
	Cold → Yes	3/4	

Attribute	Rules	Error	Total Error
Temp	Hot → No	2/4	5/14
	Mild → Yes	2/6	
	Cool → Yes	1/4	

Humidity	High	Yes 3	Decision
	Normal	Yes 6	Decision
	Normal	No 4	Decision

Attribute	Rules	Error	Total Error
Humidity	High → No	3/7	4/14
	Normal → Yes	1/7	

Windy:	False	8	Yes	6
			No	2
	True	6	Yes	3
			No	3

Attribute	Rules	Error	Total Error
Windy	True → No	3/6	5/14
	False → Yes	2/8	

Attribute	Rules	Error	Total Error
outlook	Sunny → No	2/5	4/14
	Overcast → Yes	0/4	
Temp	Rainy → Yes	2/5	5/14
	Hot → No	2/4	
Humidity	Mild → Yes	2/6	4/14
	Cool → Yes	1/4	
Windy	High → No	3/7	5/14
	Normal → Yes	1/7	
Windy	False → Yes	2/8	5/14
	True → No	3/6	

outlook

sunny

overcast

Rainy

Temp	Humidity	Windy	play	Temp	Humidity	Windy play
Hot	High	False	No	Mild	High	False Yes
Hot	High	True	No	Cool	Normal	False No
Mild	High	False	No	Cool	Normal	True No
Cool	Normal	False	Yes	Mild	Normal	False Yes
Mild	Normal	True	Yes	Mild	High	True No

Attribute	Rules	Error	Total Error
Temp	Hot \rightarrow No Mild \rightarrow No Yes Cool \rightarrow Yes	012✓ 112 011	015
Humidity	High \rightarrow No Normal \rightarrow Yes	013✓ 012✓	015
Windy	False \rightarrow No True \rightarrow Yes No	113 112	215

Attribute	Rules	Error	Total Error
Temp	Mild → Yes	1/3	2/5
	cool → Yes No	1/2	
	High → No Yes	1/2	
Humidity	Normal → Yes	1/3	2/5
	False → Yes	0/3	
	True → No	0/2	
Windy			0/5

