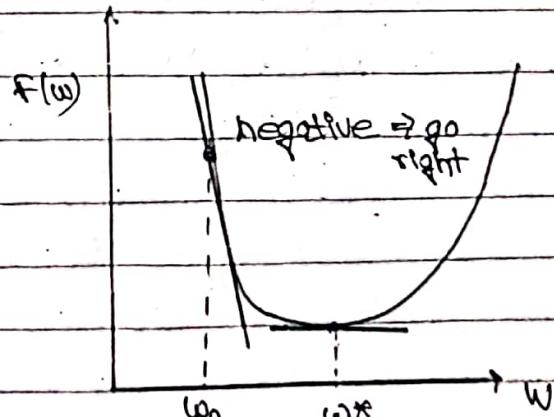
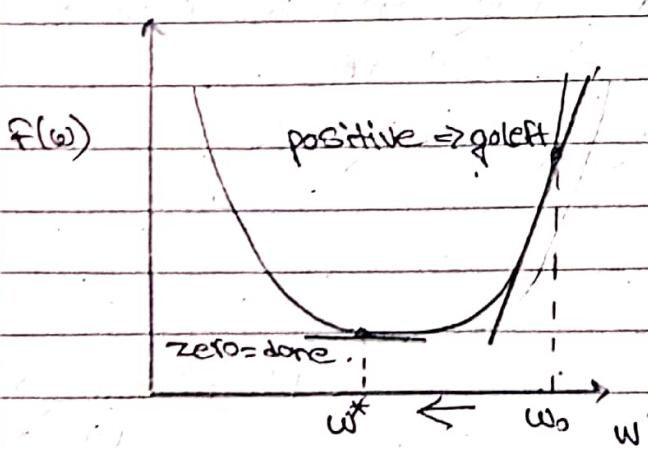
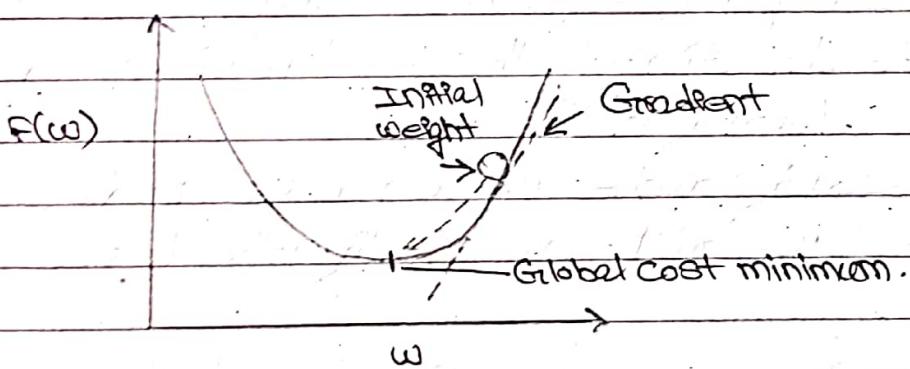


Unit: 2

Supervised Learning

* Gradient Descent:

- It is an optimization algorithm used to minimize some convex function by iteratively moving in the direction of steepest descent as defined by the negative of the gradient.
- In machine learning, we use gradient descent to update the parameters of our models. Parameters refers to coefficients in Logistic Regression & weights in neural network.



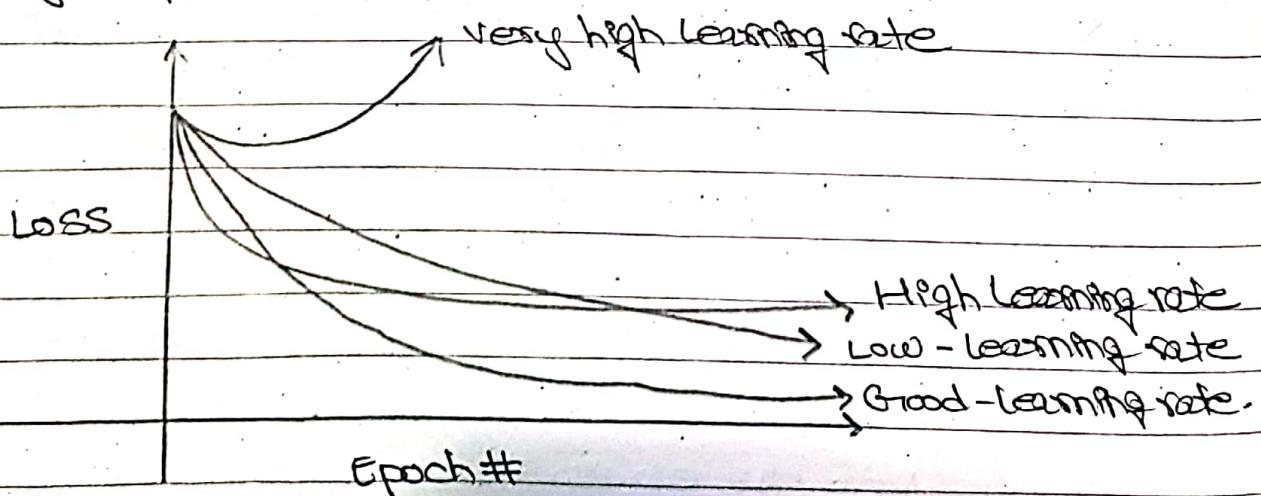
Note: We can also think of a gradient as the slope of a function.

- If f is function to be minimized (cost function), Gradient descent changes the parameters of learning model as :

$$w = w - \alpha \frac{df}{dw}$$

Where, α is learning rate & w is parameter to be optimized.

- The higher the gradient, the steeper the slope and the faster a model can learn. But if the slope is zero, the model stops learning.
- How big the steps are that Gradient descent takes into the direction of the local minimum are determined by learning rate.
- With a high learning rate we can cover more ground each step, but we risk overshooting the lowest point.
- A low learning rate is more precise, but calculating the gradient is time-consuming, so it will take us a very long time to get to the bottom.



* Gradient Descent Variations:

- Gradient Descent is an optimization algorithm which is used to find model parameters (e.g. coefficients or weights) that minimize the error of the model on the training dataset.
- It uses the error on the predictions to update the model in such a way as to reduce the error.
- It does this by making changes to the model that moves it along a gradient or slope of errors down towards a minimum error value.
- The three main flavors of gradient descent are batch, stochastic, and mini-batch.

① Batch - Gradient Descent:

- Batch - Gradient descent, computes the gradient of the cost function w.r.t. the parameters w for the entire training dataset:

$$w = w - \frac{\partial f(w, x, y)}{\partial w}$$

- As we need to calculate the gradients for the whole dataset, it performs just one update,
- Batch gradient descent is computationally fast. However, it is intractable for datasets that don't fit in memory.
- It doesn't allow us to update our model online, i.e. with new examples on-the-fly.
- Batch gradient descent is guaranteed to converge to the global minimum for convex error surfaces f to a local minimum for non-convex surfaces.

② Stochastic Gradient Descent (SGD):

- Stochastic gradient descent (SGD) in contrast performs a parameter update for each training example $x^{(i)}$ and label $y^{(i)}$. Therefore, learning happens on every example.

$$w = w - \alpha \frac{\partial f(w, x^{(i)}, y^{(i)})}{\partial w}$$

- The term stochastic indicates that one example comprising each batch is chosen at random.
- It allows us to update our model online.
- It updates the model much frequently, which is more computationally expensive than other GD variations.
- Thus, it takes significantly longer to train models on large data sets. At the same time we lose speedup due to vectorization.
- These frequent updates can result in a noisy gradient signal, which may cause the model parameters to jump around. (higher variance at training epochs)
- At the same time this behavior helps to jump to another minima.
- This ultimately complicates convergence to the exact minimum, as SGD will keep overshooting.

③ Mini-batch Gradient Descent:

- Mini-batch Gradient descent splits the training datasets into small batches that are used to calculate model error and update model coefficients.

$$w = w - \alpha \cdot \nabla f(w, x^{(i:i+n)}, y^{(i:i+n)})$$

∇w

- Mini-batch gradient descent takes the best of both Batch Gradient and SGD.
 - i) Reduces the variance of the parameter updates, which can lead to more stable convergence.
 - ii) performs less frequent parameter updates & hence is not much time consuming.
- Mini-batch gradient descent is the recommended variant of gradient descent for most applications.
- Mini-batch requires the configuration of an additional batch size Hyperparameter for the learning algorithm
 $\text{parameter} = (\text{Batch Size})$
 $\text{Hyperparameter} = (\text{Need to find experimentally})$
- Small values results in faster learning process at the cost of noise in the training process.
 (similar to batch gradient descent)
- Large values results in slow learning process with accurate estimates of the error signal.
 (similar to stochastic gradient descent)

* Linear Regression:

- Regression analysis is the process of curve fitting in which the relationship between the independent variables and dependent variables are modeled in the n^{th} degree polynomial.
- Polynomial Regression models are usually fit with the method of Least mean square (LMS) rule or Widrow-Hoff rule.
- If we assume that the relationship is a linear one and only one variable, then we can use linear equation as.

$$y = f(x) = w_0 + w_1 x \leftarrow \begin{array}{l} \text{Independent} \\ \text{variable} \end{array}$$

dependent variable

coefficients

- Let us suppose that training set contains n data points. Error function or cost function for the n data points is given by:

Let, dataset = $(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)$

$$y = f(x) = w_0 + w_1 x$$

$$y^i = w_0 + w_1 x^i \quad (\text{predicted output})$$

$$e = y^i - y = y^i - w_0 - w_1 x^i$$

where, y^i = Actual data

y = Predicted data

Error Function,

$$E = \frac{1}{2n} \sum_{i=1}^n e_i^2$$

$$E = \frac{1}{2n} \sum_{i=1}^n (y_i - w_0 - w_1 x_i)^2$$

- now, Coefficients can be determined or updated using gradient descent method as below.

$$\textcircled{1} \quad w_0 = w_0 - \alpha \frac{\partial E}{\partial w_0} = w_0 - \alpha \cdot \frac{\partial}{\partial w_0} \left(\frac{1}{2n} \sum_{i=1}^n (y_i - w_0 - w_1 x_i)^2 \right)$$

$$= w_0 + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - w_0 - w_1 x_i)$$

$$\textcircled{2} \quad w_1 = w_1 - \alpha \frac{\partial E}{\partial w_1} = w_1 + \alpha \frac{1}{n} \sum_{i=1}^n (y_i - w_0 - w_1 x_i) x_i$$

* Example: Fit a straight line through following data using SGD. Show one epoch of training.

x	1	2	3	4
f(y)	3	5	7	9

Sol

General form of linear regression equation: $y = w_0 + w_1 x$
 let us assume that initial values of parameters are:

$$w_0 = w_1 = 0$$

Iteration 1: $x = 1, y = f(x) = 3, \alpha = 0.01$

$$w_0 = w_0 + \alpha (y - w_0 - w_1 x) = 0 + 0.01 (3 - 0 - 0.1) = 0.03$$

$$w_1 = w_1 + \alpha (y - w_0 - w_1 x) x = 0 + 0.01 (3 - 0 - 0.1) \cdot 1 = 0.03$$

Iteration 2: $x = 2, y = f(x) = 5, w_0 = 0.03 \text{ and } \alpha = 0.01$

$$\begin{aligned} w_0 &= w_0 + \alpha(y - w_0 - w_0 x) \\ &= 0.03 + 0.01(5 - 0.03 - (0.03 \times 2)) \\ &= 0.0791 \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \alpha(y - w_0 - w_1 x) \\ &= 0.03 + 0.01(5 - 0.03 - (0.03 \times 2)) \times 2 \\ &= 0.1282 \end{aligned}$$

Iteration 3: $x = 3, y = f(x) = 7, w_0 = 0.0791$

$$w_1 = 0.1282$$

$$\begin{aligned} w_0 &= w_0 + \alpha(y - w_0 - w_0 x) \\ &= 0.0791 + 0.01(7 - 0.0791 \times 3) - (0.1282 \times 3) \\ &= \cancel{0.1467} \quad 0.14446 \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \alpha(y - w_0 - w_1 x) \\ &= 0.1282 + 0.01(7 - 0.0791 - (0.1282 \times 3)) \times 3 \\ &= 0.324289 \end{aligned}$$

Iteration 4: $x = 4, y = f(x) = 9, w_0 = 0.14446$

$$w_1 = 0.324289$$

$$\begin{aligned} w_0 &= w_0 + \alpha(y - w_0 - w_0 x) \\ &= 0.14446 + 0.01(9 - 0.14446 - (0.324289 \times 4)) \\ &= 0.22 \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \alpha(y - w_0 - w_1 x) \\ &= 0.324289 + 0.01(9 - 0.14446 - (0.324289 \times 4)) \times 4 \\ &= 0.63 \end{aligned}$$

* Example:

Derive weight update rule for Linear Regression for two variables.

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

cost function for linear regression,

$$E = \frac{1}{2n} \sum_{i=1}^n (y^i - w_0 - w_1 x_1 - w_2 x_2)^2$$

now, coefficient can be updated using G.D. as

$$\textcircled{1} w_0 = w_0 - \alpha \frac{\partial E}{\partial w_0} = w_0 - \alpha \cdot \frac{1}{w_0} \left(\frac{1}{2n} \sum_{i=1}^n (y^i - w_0 - w_1 x_1 - w_2 x_2)^2 \right)$$

$$= w_0 + \alpha \cdot \frac{1}{n} \sum_{i=1}^n (y^i - w_0 - w_1 x_1 - w_2 x_2)$$

$$\textcircled{2} w_1 = w_1 - \alpha \frac{\partial E}{\partial w_1} = w_1 + \alpha \cdot \frac{1}{n} \sum_{i=1}^n (y^i - w_0 - w_1 x_1 - w_2 x_2) x_1$$

$$\textcircled{3} w_2 = w_2 - \alpha \frac{\partial E}{\partial w_2} = w_2 + \alpha \cdot \frac{1}{n} \sum_{i=1}^n (y^i - w_0 - w_1 x_1 - w_2 x_2) x_2$$

* Example: Fit a straight line through the following data using SGD - show one epoch of training.

x_1	1	2	1	3
x_2	1	1	2	3
$y = f(x)$	3	5	2	5

Iteration 1: $x_1 = 1, x_2 = 1, y = f(x) = 3$

Let, $w_0 = w_1 = 0$, and $\alpha = 0.01$ then,

$$\begin{aligned} w_0 &= w_0 + \alpha (y - w_0 - w_1 x_1 - w_2 x_2) \\ &= 0 + 0.01 (3 - 0 - 0 \cdot 1 - 0 \cdot 0) \\ &= 0.03, \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \alpha (y - w_0 - w_1x_1 - w_2x_2) x_1 \\ &= 0 + 0.01(3 - 0 - 0 \cdot 1 - 0 \cdot 1) \cdot 1 \\ &= 0.03 \end{aligned}$$

$$\begin{aligned} w_2 &= w_2 + \alpha (y - w_0 - w_1x_1 - w_2x_2) x_2 \\ &= 0 + 0.01(3 - 0 - 0 \cdot 1 - 0 \cdot 1) \cdot 1 \\ &= 0.03 \end{aligned}$$

Iteration 2:, $x_1 = 2, x_2 = 1, y = 5, w_0 = 0.03$

$$w_1 = 0.03$$

$$w_2 = 0.03$$

$$\begin{aligned} w_0 &= w_0 + \alpha (y - w_0 - w_1x_1 - w_2x_2) \\ &= 0.03 + 0.01((5 - 0.03 - (0.03 \cdot 2) - (0.03 \cdot 1))) \\ &= 0.0788 \end{aligned}$$

$$\begin{aligned} w_1 &= w_1 + \alpha (y - w_0 - w_1x_1 - w_2x_2) x_1 \\ &= 0.03 + 0.01((5 - 0.03 - (0.03 \cdot 2) - (0.03 \cdot 1))) \cdot 2 \\ &= 0.1276 \end{aligned}$$

$$\begin{aligned} w_2 &= w_2 + \alpha (y - w_0 - w_1x_1 - w_2x_2) x_2 \\ &= 0.03 + 0.01((5 - 0.03 - (0.03 \cdot 2) - (0.03 \cdot 1))) \cdot 1 \\ &= 0.0788 \end{aligned}$$

Iteration 3: $x_1 = 1, x_2 = 2, y = 2, w_0 = 0.0788$

Ex

$$w_1 = 0.1276$$

$$w_2 = 0.0788$$

$$w_0 = w_0 + \alpha (y - w_0 - w_1x_1 - w_2x_2)$$

$$w_0 = 0.0788 + 0.01((2 - 0.0788 - (0.0788 \times 1) - (0.0788 \times 2))^{0.1276}) \\ = 0.09516$$

$$w_1 = w_1 + \alpha (y - w_0 - w_1 x_1 - w_2 x_2) x_1 \\ = 0.1276 + 0.01((2 - 0.0788 - (0.1276 \times 1) - (0.0788 \times 2)) \\ = 0.14396$$

$$w_2 = w_2 + \alpha (y - w_0 - w_1 x_1 - w_2 x_2) x_2 \\ = 0.0788 + 0.01(2 - 0.0788 - (0.1276 \times 1) - (0.0788 \times 2)) \times 2 \\ = 0.11152$$

Iteration 4: $x_1 = 3, x_2 = 3, y = 5, w_0 = 0.09516$

$$w_1 = 0.14396$$

$$w_2 = 0.11152$$

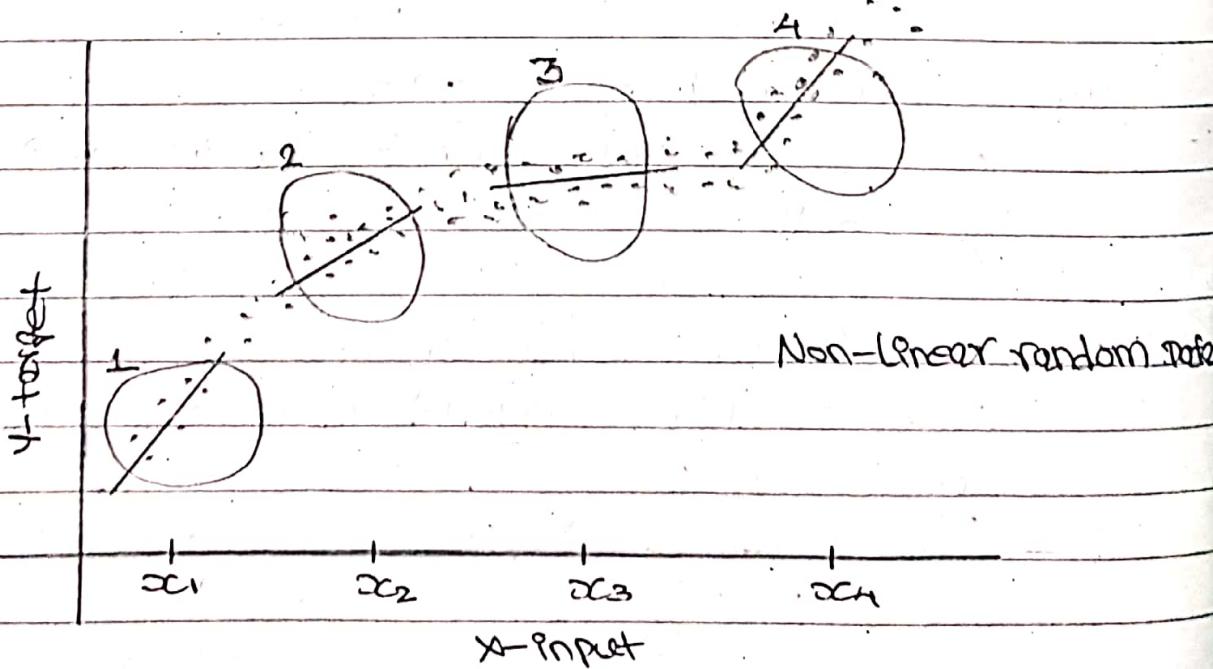
$$w_0 = w_0 + \alpha (y - w_0 - w_1 x_1 - w_2 x_2) \\ = 0.09516 + 0.01(5 - 0.09516 - (0.14396 \times 3) - (0.11152 \times 3)) \\ = 0.136544$$

$$w_1 = w_1 + \alpha (y - w_0 - w_1 x_1 - w_2 x_2) x_1 \\ = 0.14396 + 0.01(5 - 0.09516 - (0.14396 \times 3) - (0.11152 \times 3)) \times 3 \\ = 0.268112$$

$$w_2 = w_2 + \alpha (y - w_0 - w_1 x_1 - w_2 x_2) x_2 \\ = 0.11152 + 0.01(5 - 0.09516 - (0.14396 \times 3) - (0.11152 \times 3)) \times 3 \\ = 0.235672$$

* Locally weighted Regression:

- Linear regression is a supervised learning algorithm used for computing linear relationships between input (x) and output (y).
- This algorithm cannot be used for making predictions where there exists a non-linear relationship between x and y . In such case, Locally weighted Linear regression^{is} used.
- Locally weighted linear regression is a non-parametric algorithms, that is, the model does not learn a fixed set of parameters as it is done in ordinary linear regression.
- rather the parameters are computed individually for each query point x_c .
- While computing parameters a higher preference is given to the points in the training set lying in the neighbourhood of x_c than the points lying far away from x_c . i.e.



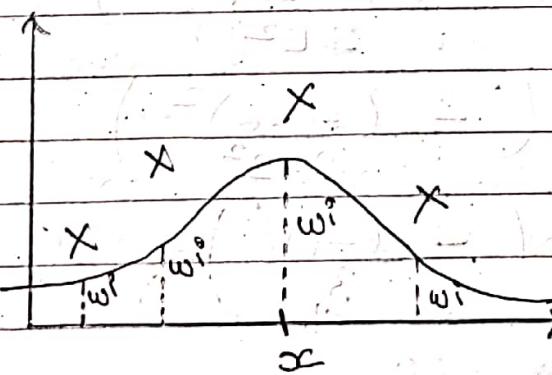
- This results in the model fitting a straight line only to the date which is near or close to the query point.
- The cost function for the locally weighted linear regression algorithm is given as below:

$$E = \frac{1}{2n} \sum_{i=1}^n w_i (y_i - \alpha_0 - \alpha_1 x_i)^2$$

- here, w_i is a non-negative weight associated with training point x_i . (The value of w_i lies between 0-1).
- For x_i lying closer to the query point the value of w_i is large, while for x_i lying away to the query point value of w_i is small. A typical choice w_i of is :

$$w_i = \exp \left(\frac{-(x_i - x_q)^2}{2\tau^2} \right)$$

where, τ called the bandwidth parameter & controls the rate at which w_i falls with distance from x .



- clearly, if $|x_i - x|$ is small w_i is close to 1 and if $|x_i - x|$ is large w_i is close to 0. Thus, the training-set-points lying closer to the query point x contribute more to the cost function than the points lying far away from x .
- Once we have cost function, we can use gradient descent algorithm to train the LWR algorithm.

* Example: Consider a query point $x = 6$ and let $x^1 = 5$, $x^2 = 4$, and $x^3 = 3$ are three points in the training set. Find the cost function for the locally weighted linear regression. (Let $T = 2$)

= Cost function for the locally weighted linear regression.

$$E = \frac{1}{2n} \sum_{i=1}^n w_i (y^i - a_0 - a_1 x^i)^2$$

Initially, we need to evaluate the value of w^i for training point $x^1 = 5$, $x^2 = 4$ and $x^3 = 3$

$$w^{(i)} = \exp \left(-\frac{(x^{(i)} - x)^2}{2T^2} \right)$$

① For $x^1 = 5$, $x = 6$, let $T = 2$

$$w^1 = \exp \left(-\frac{(x^{(1)} - x)^2}{2T^2} \right)$$

$$= \exp \left(-\frac{(5-6)^2}{2 \cdot 2^2} \right)$$

$$= \exp \left(-\frac{(-1)^2}{8} \right)$$

$$= 0.88 \text{ ..}$$

② For $x^2 = 4$, $x = 6$, let $T = 2$

$$w^2 = \exp \left(-\frac{(4-6)^2}{2 \cdot 2^2} \right)$$

$$= 0.60 \text{ ..}$$

(3) For $x^3 = 3, x = 6$, let $\tau = 2$

$$w^3 = \exp\left(-\frac{(x^{(i)} - x)^2}{2 \cdot \tau^2}\right)$$

$$= \exp\left(-\frac{(3 - 6)^2}{2 \cdot 2^2}\right)$$

$$= 0.32$$

so, the cost function is

$$E = \frac{1}{2n} \sum_{i=1}^n w_i (y^i - a_0 - a_1 x^i)^2$$

$$= \frac{1}{2n} [0.88 (y^1 - a_0 - a_1 x^1)^2 + 0.61 (y^2 - a_0 - a_1 x^2)^2 + 0.32 (y^3 - a_0 - a_1 x^3)^2]$$

Thus, the weights fall exponentially as the distance between x & $x^{(i)}$ increases

* Coefficient update of LWR:

The coefficient update rule for the locally weighted linear regression is given below:

i.e.

$$a_0 = a_0 - \frac{\partial E}{\partial a_0} = a_0 + \alpha \frac{1}{n} \sum_{i=1}^n w_i (y^i - a_0 - a_1 x^i)$$

$$a_1 = a_1 - \frac{\partial E}{\partial a_1} = a_1 + \alpha \frac{1}{n} \sum_{i=1}^n w_i (y^i - a_0 - a_1 x^i) x^i$$

* Example: Consider following dataset & show one epoch of Training using LWR.

x	3	4	5	6
y	10	15	24	

(from previous example)

= Iteration 1: $x^{(1)} = 3, w^{(1)} = 0.32, y = 10, a_0 = 0, a_1 = 0$ &
 $\alpha = 0.1$

then,

$$\begin{aligned} a_0 &= a_0 + \alpha w^i (y^i - a_0 - a_1 x^i) \\ &= 0 + 0.1 \times 0.32 (10 - 0 - 0 \times 3) \\ &= 0.32 \end{aligned}$$

$$\begin{aligned} a_1 &= a_1 + \alpha w^i (y^i - a_0 - a_1 x^i) x^i \\ &= 0 + 0.1 \times 0.32 (10 - 0 - 0) \times 3 \\ &= 0.96 \end{aligned}$$

Iteration 2: $x^{(2)} = 4, w^{(2)} = 0.60, y = 15, a_0 = 0.32$

$$a_1 = 0.96$$

$$\alpha = 0.1$$

$$\begin{aligned} a_0 &= a_0 + \alpha w^i (y^i - a_0 - a_1 x^i) \\ &= 0.32 + 0.1 \times 0.60 (15 - 0.32 - 0.96 \times 4) \\ &= 0.97 \end{aligned}$$

$$\begin{aligned} a_1 &= a_1 + \alpha w^i (y^i - a_0 - a_1 x^i) x^i \\ &= 0.96 + 0.1 \times 0.60 (15 - 0.32 - 0.96 \times 4) \times 4 \\ &= 3.60 \end{aligned}$$

Iteration 3^o: $x^{(i)} = 5$, $w^{(i)} = 0.88$, $y = 24$, $a_0 = 0.97$

$$\alpha_1 = 3.60$$

$$\alpha = 0.1$$

$$a_0 = a_0 + \alpha w_i (y^i - a_0 - \alpha_1 x^i)$$

$$= 0.97 + 0.1 \times 0.88 (24 - 0.97 - 3.60 \times 5)$$

$$= 1.41$$

$$\alpha_1 = \alpha_1 + \alpha w_i (y^i - a_0 - \alpha_1 x^i) x^i$$

$$= 3.60 + 0.1 \times 0.88 (24 - 0.97 - 3.60 \times 5) \times 5$$

$$= 5.81$$

* Logistic Regression:

- Logistic regression is one of the most popular machine learning algorithms for binary classification. This is because it is a simple algorithm that performs very well on a wide range of problems.
- We want to predict a variable $y \in \{0, 1\}$, where 0 is called negative class, while 1 is called positive class. Such task is known as binary classification.
- The heart of the logistic regression technique is logistic function and is defined as:

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

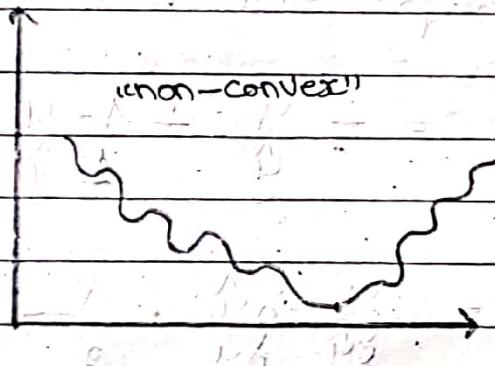
- Logistic function transforms the input into the range [0, 1]. Smallest negative numbers results in values close to zero and the larger positive numbers results in values close to one.
- If there are two input variables, logistic regression has two coefficients just like linear regression. i.e.

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

- unlike linear regression, the output is transformed into a probability using the logistic function!

$$\hat{y} = \sigma(y) = \frac{1}{1+e^{-y}}$$

- If the probability is > 0.5 we can take the output as a prediction for the class 1, otherwise the prediction is for the class 0.
- The job of the learning algorithm will be to discover the best values for the coefficients (w_0, w_1 , and w_2) based on the training data.
- Unfortunately, we can't use the cost function MSE in logistic regression. It's because our prediction function is non-linear (due to sigmoid transform).
- Squaring this prediction as we do in MSE results in a non-convex function with many local minimums. If our cost function has many local minimums, gradient descent may not find the optimal global minimum.



- Given the training set $S (x^{(1)}, y^{(1)}, (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$ we want $\hat{y}^{(i)} \approx y^{(i)}$
- For logistic regression, we use the following loss function or error function

$$L(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y=1 \\ -\log(1-\hat{y}) & \text{if } y=0 \end{cases}$$

- In case $y=1$, the output (the cost to pay) approaches to 0

as \hat{y} approaches to 1. Conversely, the cost to pay goes to infinity as \hat{y} approaches to 0.

- This is a desirable property: we want a bigger penalty as the algorithm predicts something far away from the actual value. The same intuition applies when $y=0$.
- Thus, Cost function for logistic regression is given as:

$$L(y, \hat{y}) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

* Logistic Regression Gradients/derivatives:

$$\frac{\partial L(y, \hat{y})}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

$$\frac{\partial L(y, \hat{y})}{\partial y} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} = \left\{ -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right\} \left\{ \hat{y}(1-\hat{y}) \right\}$$

$$= \hat{y} - y$$

if $g(x)$ is logistic function $g'(x) = g(x)(1-g(x))$

$$\frac{\partial L}{\partial w_0} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} \frac{\partial y}{\partial w_0} = (\hat{y} - y_0)$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial y} \frac{\partial y}{\partial w_i} = \alpha_i(\hat{y} - y)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial y} \cdot \frac{\partial y}{\partial w_2} = x_2 (\hat{y} - y)$$

Thus, parameters of logistic regression are updated as below:

$$w_0 = w_0 - \alpha \frac{\partial L}{\partial w_0} = w_0 - \alpha (\hat{y} - y)$$

$$w_1 = w_1 - \alpha \frac{\partial L}{\partial w_1} = w_1 - \alpha (\hat{y} - y) x_1$$

$$w_2 = w_2 - \alpha \frac{\partial L}{\partial w_2} = w_2 - \alpha (\hat{y} - y) x_2$$

* Example: Fit the logistic regression model through the following data. Show one epoch of training.

x_1	x_2	class(y)
0.78	0.69	1
0.67	1.00	1
0.00	0.00	-1
0.22	0.14	-1

= Sol

General form of logistic regression equation is

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

let us assume that initial values of parameters are:

$$w_0 = w_1 = w_2 = 0$$

Iteration 1: $x_1 = 0.78, x_2 = 0.69, y = 1, \alpha = 0.1$

$$y = w_0 + w_1 x_1 + w_2 x_2 = 0 + 0 \times 0.78 + 0 \times 0.69 = 0$$

So,

$$\hat{y} = \frac{1}{1 + e^{-y}} = \frac{1}{1 + e^{-0}} = 0.5$$

$$\begin{aligned} w_0 &= w_0 - \alpha(\hat{y} - y) = 0 - 0.1(0.5 - 1) = 0.05 \\ w_1 &= w_1 - \alpha(\hat{y} - y)x_1 = 0 - 0.1(0.5 - 1) * 0.78 = 0.039 \\ w_2 &= w_2 - \alpha(\hat{y} - y)x_2 = 0 - 0.1(0.5 - 1) * 0.69 = 0.0345 \end{aligned}$$

Iteration 2:

$$x_1 = 0.67, x_2 = 1, y = 1, \alpha = 0.1, w_0 = 0.05$$

$$w_1 = 0.039$$

$$w_2 = 0.0345$$

$$\begin{aligned} y &= w_0 + w_1 x_1 + w_2 x_2 \\ &= 0.05 + (0.039 * 0.67) + (0.0345 * 1) \\ &= 0.11063 \end{aligned}$$

$$\hat{y} = \frac{1}{1+e^{-y}} = 0.5276$$

$$\begin{aligned} w_0 &= w_0 - \alpha(\hat{y} - y) = 0.05 - 0.1(0.5276 - 1) = 0.097 \\ w_1 &= w_1 - \alpha(\hat{y} - y)x_1 = 0.039 - 0.1(0.5276 - 1) * 0.67 \\ &= 0.0706 \end{aligned}$$

$$\begin{aligned} w_2 &= w_2 - \alpha(\hat{y} - y)x_2 = 0.0345 - 0.1(0.5276 - 1)x_1 \\ &= 0.08174 \end{aligned}$$

Iteration 3: $x_1 = 0, x_2 = 0, y = 0, \alpha = 0.1, w_0 = 0.097$

$$w_1 = 0.0706$$

$$w_2 = 0.08174$$

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

$$\begin{aligned} &= 0.097 + (0.0706 * 0) + (0.08174 * 0) \\ &= 0.097 \end{aligned}$$

$$\hat{y} = \frac{1}{1+e^{-y}} = 0.5243$$

$$w_0 = w_0 - \alpha (\hat{y} - y) = 0.097 - 0.1 (0.5243 - 0) = 0.0445$$

$$w_1 = w_1 - \alpha (\hat{y} - y)x_1 = 0.0706 - 0.1 (0.5243 - 0) \times 0.22 = 0.0706$$

$$w_2 = w_2 - \alpha (\hat{y} - y)x_2 = 0.08174 - 0.1 (0.5243 - 0) \times 0.14 = 0.08174$$

Iteration 4:

$$x_1 = 0.22, x_2 = 0.14, y = 0, \alpha = 0.1, w_0 = 0.0445$$

$$w_1 = 0.0706$$

$$w_2 = 0.08174$$

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

$$= 0.0445 + (0.0706 \times 0.22) + (0.08174 \times 0.14)$$

$$= 0.0715$$

$$\hat{y} = \frac{1}{1+e^{-y}} = 0.5178$$

$$w_0 = w_0 - \alpha (\hat{y} - y) = 0.0445 - 0.1 (0.5178 - 0) = -0.072$$

$$w_1 = w_1 - \alpha (\hat{y} - y)x_1 = 0.0706 - 0.1 (0.5178 - 0) \times 0.22 = 0.0592$$

$$w_2 = w_2 - \alpha (\hat{y} - y)x_2 = 0.08174 - 0.1 (0.5178 - 0) \times 0.14 = 0.0744$$

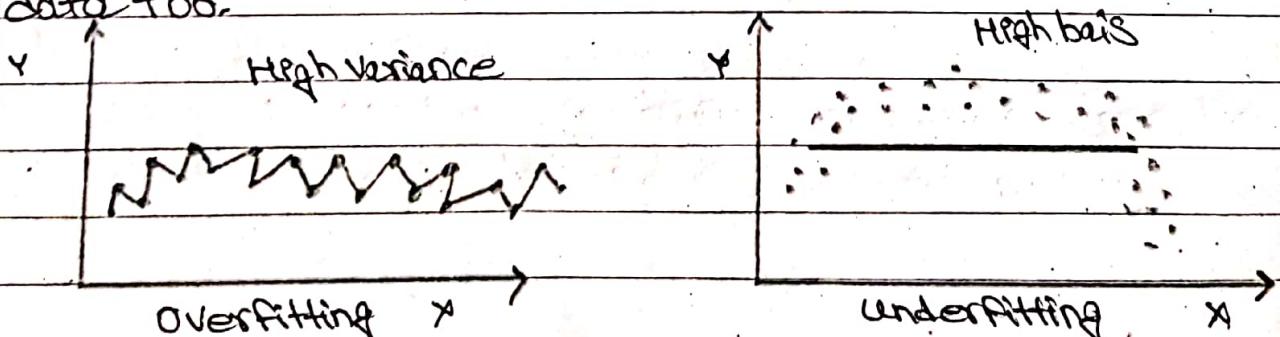
* Overfitting and Underfitting:

① Overfitting:

- overfitting is the result of using an excessively complicated model.
- It happens when the model learns the details and noise of training data to the extent that it negatively impacts the performance of the model on new data.
- This means that the model learns noise or random fluctuations in the training data as concepts.
- This negatively impacts the model's ability to generalize because these concepts do not apply to new data.

② Underfitting:

- Underfitting is the result of using an excessively simple model or using very few training samples.
- In such situations, a machine learning algorithm cannot capture the underlying trend of the data.
- Thus, it refers to a model that can neither model the training data nor generalize to new data.
- Obviously, an underfitted machine learning model is not a suitable model for making predictions because it has poor performance on the training data too.



* Performance Measures of classification:

- Before selecting a classification algorithm to solve a problem, we need to evaluate its performance. Confusion matrix is widely used for computing performance measures of classification algorithm.

* Confusion Matrix:

- A Confusion matrix is an $N \times N$ matrix used for evaluating the performance of a classification model, where N is the number of target classes.
- The matrix compares the actual target values with those predicted by the machine learning model. This gives us a holistic view of how well our classification model is performing and what kinds of errors it is making.
- For a binary classification problem, we would have a 2×2 matrix as shown below with 4 values:

		Actual classes	
		Positive	Negative
predicted classes	positive	TP	FP
	Negative	FN	TN

① True Positive (TP):

It represents correctly classified positive classes. Both actual and predicted class are positive here.

② False Positive (FP):

It represents incorrectly classified positive classes. These are the positive classes predicted by the model that were actually negative. This is called Type I error.

③ True Negative (TN):

It represents correctly classified Negative classes. Both actual & predicted class are negative here.

④ False Negative (FN):

It represents incorrectly classified negative classes. These are the negative classes predicted by the model that were actually positive. This is called Type II error.

* Performance measures of classifications:

Four widely used performance measures used for evaluating classification models are:

- ① Accuracy,
- ② Recall,
- ③ Precision,
- ④ F1-Score

① Accuracy:

It is the percentage of correct predictions made by the model and is given as below:

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

Total no. of instances

② Precision:

It is percentage of predicted positives that are actually positive and is given by:

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

Total number of predicted positive

③ Recall:

It is the percentage of actual positives that are correctly classified by the model and is given as below:

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

Total number of actual positives.

④ F1-Score:

It is the harmonic mean of recall and precision. It becomes high only when both precision and recall are high. This score is given by:

$$F1 = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$$

* Example: Suppose that we have to classify 100 people as pregnant or not pregnant. This includes 40 pregnant women & the remaining 60 are not pregnant. Out of 40 pregnant women 30 pregnant women are classified correctly and the remaining 10 percent women are classified as not pregnant by the machine learning algorithm. On the other hand, out of 60 people in the not pregnant category, 55 are classified as not pregnant & the remaining 5 are classified pregnant. Compute accuracy, precision, recall and F1-Score for above example.

= Sol

Let, the pregnant women belongs to Positive class and non-pregnant women belongs to Negative class then,

$$\text{True positive} = 30$$

$$\text{False positive} = 5$$

$$\text{True Negative} = 55$$

$$\text{False Negative} = 10$$

		Actual class	
		Pregnant	Non-pregnant
predicted class	Pregnant	30 (T.P)	5 (F.P)
	Non-pregnant	10 (F.N)	55 (T.N)

$$\textcircled{1} \quad \text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{30 + 55}{30 + 5 + 10 + 10} = \frac{85}{100} = 0.85$$

$$\textcircled{2} \quad \text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{30}{30 + 5} = \frac{30}{35} = 0.857$$

$$\textcircled{3} \quad \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{30}{30 + 10} = \frac{30}{40} = 0.75$$

$$\textcircled{4} \quad \text{F1-score} = \frac{2 \times \text{Recall} \times \text{Precision}}{\text{Recall} + \text{precision}} = \frac{2 \times 0.75 \times 0.857}{0.75 + 0.857} = 0.799$$

Ans: 0.799

* Matrix Derivatives:

- Let A be a $m \times n$ matrix and $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a mapping function from m -by- n matrices to a real number. i.e.

$$A = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & A_{22} & \dots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \dots & A_{mn} \end{bmatrix}$$

- Now, Gradient of f with respect to A is also a $m \times n$ matrix and is given as below:

$$\nabla_A f_A = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} & \dots & \frac{\partial f}{\partial A_{1n}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} & \dots & \frac{\partial f}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{m1}} & \frac{\partial f}{\partial A_{m2}} & \dots & \frac{\partial f}{\partial A_{mn}} \end{bmatrix}$$

* Example: Consider the following 2×2 matrix,

$$\begin{bmatrix} w & x \\ y & z \end{bmatrix}$$

- Let $f(w, x, y, z) = 2w + 3x + yz$

- $\nabla_A f_A = \begin{bmatrix} 2 & 3 \\ z & y \end{bmatrix}$

* Trace of Matrix:

- For a $n \times n$ matrix A , trace of A is defined as sum of diagonal elements
i.e. $\text{tr. } A = \sum_{i=1}^n a_{ii}$
- For a real number a $\text{tr. } a = a \cdot 1$
- If multiplication of $A \& B$ is square matrix then $\text{tr. } AB = \text{tr. } BA$

Verify that $\text{tr. } AB = \text{tr. } BA$

Let, A & B are matrices of order 2×3 & 3×2 respectively,

i.e. $A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}$

$$B = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix}$$

$$AB = \begin{bmatrix} a_{11} \cdot b_{11} + a_{12} \cdot b_{21} + a_{13} \cdot b_{31} & a_{11} \cdot b_{12} + a_{12} \cdot b_{22} + a_{13} \cdot b_{32} \\ a_{21} \cdot b_{11} + a_{22} \cdot b_{21} + a_{23} \cdot b_{31} & a_{21} \cdot b_{12} + a_{22} \cdot b_{22} + a_{23} \cdot b_{32} \end{bmatrix}$$

$$\text{tr. } AB = a_{11} \cdot b_{11} + a_{12} \cdot b_{21} + a_{13} \cdot b_{31} + a_{21} \cdot b_{12} + a_{22} \cdot b_{22} + a_{23} \cdot b_{32}$$

Similarly

$$BA = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}$$

$$= \begin{bmatrix} b_{11} \cdot a_{11} + b_{12} \cdot a_{21} & b_{11} \cdot a_{12} + b_{12} \cdot a_{22} & b_{11} \cdot a_{13} + b_{12} \cdot a_{23} \\ b_{21} \cdot a_{11} + b_{22} \cdot a_{21} & b_{21} \cdot a_{12} + b_{22} \cdot a_{22} & b_{21} \cdot a_{13} + b_{22} \cdot a_{23} \\ b_{31} \cdot a_{11} + b_{32} \cdot a_{21} & b_{31} \cdot a_{12} + b_{32} \cdot a_{22} & b_{31} \cdot a_{13} + b_{32} \cdot a_{23} \end{bmatrix}$$

$$\text{tr. } BA = [b_{11} \cdot a_{11} + b_{12} \cdot a_{21} + b_{21} \cdot a_{12} + b_{22} \cdot a_{22} + b_{31} \cdot a_{13} + b_{32} \cdot a_{23}]$$

$$= [a_{11} \cdot b_{11} + a_{12} \cdot b_{21} + a_{13} \cdot b_{31} + a_{21} \cdot b_{12} + a_{22} \cdot b_{22} + a_{23} \cdot b_{32}]$$

$$\therefore \text{tr. } AB = \text{tr. } BA$$

* Corollaries of $\text{tr. } AB = \text{tr. } BA$

$$\text{tr. } ABC = \text{tr. } CAB = \text{tr. } BCA$$

$$\text{tr. } ABCD = \text{tr. } DABC = \text{tr. } BCDA$$

* Prove that $\text{tr. } ABC = \text{tr. } CAB = \text{tr. } BCA$

i.e.,

$$\begin{aligned}\text{tr. } ABC &= \text{tr. } (AB)C \\ &= \text{tr. } CAB\end{aligned}$$

$$\text{Similarly, } \text{tr. } ABC = \text{tr. } A(BC)$$

$$\stackrel{A}{\bullet} \stackrel{B}{\bullet}$$

$$= \text{tr. } BCA$$

Thus, $\text{tr. } ABC = \text{tr. } CAB = \text{tr. } BCA$ proved.

* Prove that $\text{tr. } ABCD = \text{tr. } DABC = \text{tr. } BCDA$

i.e.,

$$\text{tr. } ABCD = \text{tr. } (ABC)D$$

$$= \text{tr. } DABC$$

$$\text{tr. } ABCD = \text{tr. } A(BCD)$$

$$= \text{tr. } BCDA$$

Thus, $\text{tr. } ABCD = \text{tr. } DABC = \text{tr. } BCDA$ proved.

* If A and B are square matrices & a is a real number, then trace operator exhibits following properties.

$$\text{i.e., } \text{① } \text{tr} \cdot A = \text{tr} \cdot A^T$$

$$\text{② } \text{tr} \cdot (A+B) = \text{tr} \cdot A + \text{tr} \cdot B$$

$$\text{③ } \text{tr} \cdot a \cdot A = a \cdot \text{tr} \cdot A$$

$$\text{① } \text{tr} \cdot A = \text{tr} \cdot A^T$$

Sol

$$\text{let, } A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

$$\text{tr} \cdot A = \sum_{i=1}^n a_{ii}$$

Again,

$$A^T = \begin{bmatrix} a_{11} & a_{21} & \dots & a_{n1} \\ a_{12} & a_{22} & \dots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

$$\text{tr} \cdot A^T = \sum_{i=1}^n a_{ii}$$

$$\text{Thus, } \boxed{\text{tr} \cdot A = \text{tr} \cdot A^T} \quad \text{i.e., } \boxed{\sum_{i=1}^n a_{ii}}$$

$$\textcircled{2} \quad \text{tr}(A+B) = \text{tr}. A + \text{tr}. B$$

= let,

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & & & \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}$$

$$B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & & & \\ b_{n1} & b_{n2} & \dots & b_{nn} \end{bmatrix}$$

now,

$$A+B = \begin{bmatrix} a_{11}+b_{11} & a_{12}+b_{12} & \dots & a_{1n}+b_{1n} \\ a_{21}+b_{21} & a_{22}+b_{22} & \dots & a_{2n}+b_{2n} \\ \vdots & & & \\ a_{n1}+b_{n1} & a_{n2}+b_{n2} & \dots & a_{nn}+b_{nn} \end{bmatrix}$$

$$\text{tr}(A+B) = [a_{11}+b_{11} + a_{22}+b_{22} + \dots + a_{nn}+b_{nn}]$$

$$\text{tr}. (A+B) = \sum_{i=1}^n a_{ii} + b_{ii}$$

Then,

$$\text{tr}. A = a_{11} + a_{22} + \dots + a_{nn} = \sum_{i=1}^n a_{ii}$$

$$\text{tr}. B = b_{11} + b_{22} + \dots + b_{nn} = \sum_{i=1}^n b_{ii}$$

$$\therefore \text{tr}. A + \text{tr}. B = \sum_{i=1}^n a_{ii} + \sum_{i=1}^n b_{ii}$$

$$= \sum_{i=1}^n a_{ii} + b_{ii}$$

$$\text{Thus, } [\text{tr}. (A+B) = \text{tr}. A + \text{tr}. B]$$

$$\textcircled{3} \quad \text{tr. } aA = a\text{tr. } A$$

= Let,

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$a \cdot A = \begin{bmatrix} a \cdot a_{11} & a \cdot a_{12} \\ a \cdot a_{21} & a \cdot a_{22} \end{bmatrix}$$

$$\text{tr. } aA = a \cdot a_{11} + a \cdot a_{22}$$

now,

$$\text{tr. } A = a_{11} + a_{22}$$

thus,

$$a \cdot \text{tr. } A = a \cdot a_{11} + a \cdot a_{22}$$

* Some facts about matrix derivatives:

$$\textcircled{1} \quad \nabla_A \text{tr} AB = B^T$$

$$\textcircled{2} \quad \nabla_A^T f(A) = (\nabla_A f(A))^T$$

$$\textcircled{3} \quad \nabla_A \text{tr} ABC^T = CAB + C^T AB^T$$

$$\textcircled{4} \quad \nabla_A |A| = |A|(A^{-1})^T$$

For more details refer

① Prove that $\nabla_A \text{tr} \cdot AB = B^T$

= let order of A is 3×2 & order of B is 2×3 .

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix}, \quad B = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \end{bmatrix}$$

$$AB = \begin{bmatrix} a_{11} \cdot b_{11} + a_{12} \cdot b_{21} & a_{11} \cdot b_{12} + a_{12} \cdot b_{22} & a_{11} \cdot b_{13} + a_{12} \cdot b_{23} \\ a_{21} \cdot b_{11} + a_{22} \cdot b_{21} & a_{21} \cdot b_{12} + a_{22} \cdot b_{22} & a_{21} \cdot b_{13} + a_{22} \cdot b_{23} \\ a_{31} \cdot b_{11} + a_{32} \cdot b_{21} & a_{31} \cdot b_{12} + a_{32} \cdot b_{22} & a_{31} \cdot b_{13} + a_{32} \cdot b_{23} \end{bmatrix}$$

$$\text{tr} \cdot AB = a_{11} \cdot b_{11} + a_{12} \cdot b_{21} + a_{21} \cdot b_{12} + a_{22} \cdot b_{22} + a_{31} \cdot b_{13} + a_{32} \cdot b_{23}$$

now,

$$\nabla_A \text{tr} \cdot AB = \begin{bmatrix} b_{11} & b_{21} \\ b_{12} & b_{22} \\ b_{13} & b_{23} \end{bmatrix}$$

$$B^T = \begin{bmatrix} b_{11} & b_{21} \\ b_{12} & b_{22} \\ b_{13} & b_{23} \end{bmatrix}$$

Thus $\nabla_A \text{tr} \cdot AB = B^T$ proved.

* Deriving Normal Equation For Linear Regression:

= Given, a training set, let x be the $m \times n$ matrix that contains input values of training examples in its rows and let \vec{y} be a m -dimensional vector containing target values.

$$\text{Given } X = \begin{bmatrix} (x^1)^T & y^1 \\ (x^2)^T & y^2 \\ \vdots & \vdots \\ (x^m)^T & y^m \end{bmatrix}$$

We know that,

$$y = f(x) = (x^i)^T w, \text{ where } w \text{ is coefficient vector.}$$

now,

$$\begin{aligned} \text{error} &= f(x) - \vec{y} \\ &= x^i w - \vec{y} \end{aligned}$$

$$= \begin{bmatrix} (x^1)^T w & y^1 \\ (x^2)^T w & y^2 \\ \vdots & \vdots \\ (x^m)^T w & y^m \end{bmatrix}$$

now, Loss function can be written as :

$$L = \frac{1}{2} \sum_{i=1}^m (f(x_i) - y_i)^2$$

$$= \frac{1}{2} (x\omega - \vec{y})^T (x\omega - \vec{y})$$

to minimize loss, we need to calculate derivatives / gradient of L w.r.t. ω .

$$\nabla_\omega L = \frac{1}{2} \nabla_\omega (x\omega - \vec{y})^T (x\omega - \vec{y})$$

$$= \frac{1}{2} \nabla_\omega (\omega^T x^T - \vec{y}^T) (x\omega - \vec{y}) \quad [\because (AB)^T = B^T A]$$

$$= \frac{1}{2} \nabla_\omega (\omega^T x^T x\omega - \omega^T x^T \vec{y} - \vec{y}^T x\omega + \vec{y}^T \vec{y})$$

$$= \frac{1}{2} \nabla_\omega \text{tr}(\omega^T x^T x\omega - \omega^T x^T \vec{y} - \vec{y}^T x\omega + \vec{y}^T \vec{y})$$

$$= \frac{1}{2} \nabla_\omega (\text{tr} \omega^T x^T x\omega - 2 \text{tr} \vec{y}^T x\omega) \quad [\because \nabla_\omega \vec{y}^T \vec{y} = 0]$$

$$\text{tr } A = \text{tr } A^T$$

$$\text{tr} (\omega^T x^T \vec{y}) = \text{tr} (\omega^T x^T \vec{y})^T$$

$$= \text{tr} (\vec{y}^T x\omega)$$

$$= \frac{1}{2} (x^T x\omega + x^T x\omega - 2 x^T \vec{y}) \quad [\because \nabla_A \text{tr } AB A^T C = B^T A^T C^T + B A^T C]$$

$$= x^T x\omega - x^T \vec{y}$$

At minima, gradient / derivatives is zero :

$$\text{thus, } x^T x\omega - x^T \vec{y} = 0$$

$$\Rightarrow \omega = \frac{x^T \vec{y}}{x^T x}$$

$$= (x^T x)^{-1} \cdot x^T \vec{y}$$

→ This is called normal eqn for linear regression.

Note: We can use above equation to find values of coefficients.

Ex: Example:

x_1	1	2	1	3
x_2	1	1	2	3
y	3	5	2	5

Use normal equation to find coefficient of linear regression equation fitted through above example.

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 1 & 2 \\ 3 & 3 \end{bmatrix}, \vec{y} = \begin{bmatrix} 3 \\ 5 \\ 2 \\ 5 \end{bmatrix}$$

Augment \vec{x} for w_0 i.e., $y = w_0 + w_1x_1 + w_2x_2$

Padding for w_0 value.

$$\vec{x} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \\ 1 & 3 & 3 \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} 3 \\ 5 \\ 2 \\ 5 \end{bmatrix}$$

We have, normal equation i.e.

$$w = (x^T x)^{-1} x^T \vec{y}$$

$$= \left(\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 1 & 3 \\ 1 & 1 & 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \\ 1 & 3 & 3 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \\ 1 & 3 & 3 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 5 \\ 2 \\ 5 \end{bmatrix}$$

$$= \frac{1+1+1+1}{1+2+1+3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \\ 1 & 3 & 3 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 5 \\ 2 \\ 5 \end{bmatrix}$$

$$\frac{1+2+1+3}{1+4+1+9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \\ 1 & 3 & 3 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 5 \\ 2 \\ 5 \end{bmatrix}$$

$$\frac{1+1+2+3}{1+2+2+9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \\ 1 & 3 & 3 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 5 \\ 2 \\ 5 \end{bmatrix}$$

$$= \begin{bmatrix} 4 & 7 & 7 \\ 7 & 15 & 14 \\ 7 & 14 & 15 \end{bmatrix}^{-1} \begin{bmatrix} 3+5+2+5 \\ 3+10+2+15 \\ 3+5+4+15 \end{bmatrix}$$

$$= \begin{bmatrix} 4 & 7 & 7 \\ 7 & 15 & 14 \\ 7 & 14 & 15 \end{bmatrix}^{-1} \begin{bmatrix} 15 \\ 30 \\ 27 \end{bmatrix}$$

Note,

$$\begin{bmatrix} 4 & 7 & 7 \\ 7 & 15 & 14 \\ 7 & 14 & 15 \end{bmatrix}^{-1} = \frac{1}{\det(A)} \text{Adj.}(A)$$

$$\det(A) = 4 \begin{vmatrix} 15 & 14 \\ 14 & 15 \end{vmatrix} - 7 \begin{vmatrix} 7 & 14 \\ 7 & 15 \end{vmatrix}$$

$$+ 7 \begin{bmatrix} 7 & 15 \\ 7 & 14 \end{bmatrix}$$

$$= 4(225 - 106) - 7(105 - 98) + 7(98 - 105)$$

$$= 116 - 49 - 49$$

= 18 ≠ 0 then, matrix has inverse.

now,

$$\text{Adj} A =$$

$$\begin{bmatrix} 15 & 14 \\ 14 & 15 \end{bmatrix} \cdot \begin{bmatrix} 7 & 14 \\ 7 & 15 \end{bmatrix} = \begin{bmatrix} 7 & 15 \\ 7 & 14 \end{bmatrix}$$

$$+ = +29 \quad - = -7 \quad + = -7$$

$$\begin{bmatrix} 7 & 7 \\ 14 & 15 \end{bmatrix} \cdot \begin{bmatrix} 4 & 7 \\ 7 & 15 \end{bmatrix} = \begin{bmatrix} 4 & 7 \\ 7 & 14 \end{bmatrix}$$

$$- = -7 \quad + = 11 \quad + = -7$$

$$\begin{bmatrix} 7 & 7 \\ 15 & 14 \end{bmatrix} \cdot \begin{bmatrix} 4 & 7 \\ 7 & 14 \end{bmatrix} = \begin{bmatrix} 4 & 7 \\ 7 & 15 \end{bmatrix}$$

$$+ = -7 \quad - = -7 \quad + = 11$$

$$\text{Adj} A = \begin{bmatrix} 29 & -7 & -7 \\ -7 & 11 & -7 \\ -7 & -7 & 11 \end{bmatrix}^T = \begin{bmatrix} 29 & -7 & -7 \\ -7 & 11 & -7 \\ -7 & -7 & 11 \end{bmatrix}$$

now, $A^{-1} = \frac{1}{|A|} \text{adj. } A$

$$R^2 A^{-1} = \frac{1}{|A|} \text{adj. } A$$

$$= \frac{1}{18}$$

$$\begin{bmatrix} 29 & -7 & -7 \\ -7 & 11 & -7 \\ -7 & -7 & 11 \end{bmatrix}$$

$$= \frac{29}{18} \quad \frac{-7}{18} \quad \frac{-7}{18}$$

$$= \frac{-7}{18} \quad \frac{15}{18} \quad \frac{-7}{18}$$

$$= \frac{-7}{18} \quad \frac{-7}{18} \quad \frac{15}{18}$$

$$= \begin{bmatrix} 1.61 & -0.38 & -0.38 \\ -0.38 & 0.61 & -0.38 \\ -0.38 & -0.38 & 0.61 \end{bmatrix} \quad \begin{matrix} 15 \\ 30 \\ 27 \end{matrix}$$

$$W = \begin{bmatrix} 0.49 \\ 0.34 \\ -0.63 \end{bmatrix}$$

* Parametric vs Non-Parametric Algorithms:

① Parametric Algorithms:

- General form of learning model is assumed or General form of mapping function is assumed.
- Number of parameters are fixed, we just need to adjust values of parameters.
- number of parameters is independent of size of training data.
- fast to train and predict compared to non-parametric
- can perform satisfactorily with small data size,
- example:- Linear regression, Locally weighted logistic regression, i.e $y = w_0 + w_1x_1 + w_2x_2$

② Non-parametric Algorithm:

- No assumption about mapping function
- parameters are not only adjustable rather size of the parameters may change.
- slow compared to parametric algorithms.
- Need a large volume of training data to learn.
- example:- KNN (K-nearest neighbor) algorithm or classifier.

* Probabilistic Implementation of Least-Square Regression

In case of least square regression, Input and output variables are related via following equations.

$$y^i = \omega^T x^i + \epsilon^i$$

Where, ϵ^i is the error term that captures unmodeled effect & random noise.

Let us assume that ϵ^i follows IID (Independently and Identically distributed). Thus, according to Gaussian distributed distribution with zero mean & σ^2 standard deviation, probability distribution function for ϵ^i can be written as below:

$$P(\epsilon^i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\epsilon^i)^2}{2\sigma^2}\right)$$

$$\Rightarrow P(y^i/x^i, \omega) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^i - \omega^T x^i)^2}{2\sigma^2}\right)$$

If we view above function as function of ω , it is called likelihood function and is given by:

$$L(\omega) = P(\vec{y} / \vec{x}, \omega)$$

$$= \prod_{i=1}^m P(y^i/x^i, \omega)$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^i - \omega^T x^i)^2}{2\sigma^2}\right)$$

Taking log on both sides,

$$\log(L(\omega)) = \log \left(\prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y_i - \omega^T x_i)^2}{2\sigma^2} \right) \right)$$

$$\Rightarrow L(\omega) = \prod_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y_i - \omega^T x_i)^2}{2\sigma^2} \right) \right)$$

a**b**

$$L(\omega) = \log(L(\omega))$$

$$= m \cdot \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \omega^T x_i)^2$$

$$= m \cdot \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y_i - \omega^T x_i)^2$$

Thus, maximizing above equation is equivalent to minimizing the term $\frac{1}{2} \sum_{i=1}^m (y_i - \omega^T x_i)^2$

i.e., we can say that maximizing likelihood is equivalent to minimizing Squared Sum of error (i.e. least square method).

* Perception learning Algorithm:

- Perception is a single neuron that is used for binary classification.
- It can only classify linearly separable patterns.
- Perception rule can be derived from logistic regression as below.
- In logistic regression \hat{y} is predicted as below:

$$\hat{y} = g(y) = \frac{1}{1+e^{-y}}$$

where, $y = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$

- If we use following definition of g then, it will be perception learning rule,

$$g(y) = \begin{cases} 1 & y \geq 0 \\ 0 & y < 0 \end{cases}$$

- Thus, weight update rule for logistic regression is same as weight update rule for perception, which is given as below:

$$w_0 = w_0 + \alpha \frac{\partial l}{\partial w_0}$$

$$w_1 = w_1 + \alpha \frac{\partial l}{\partial w_1}$$

$$w_2 = w_2 + \alpha \frac{\partial l}{\partial w_2}$$

where,

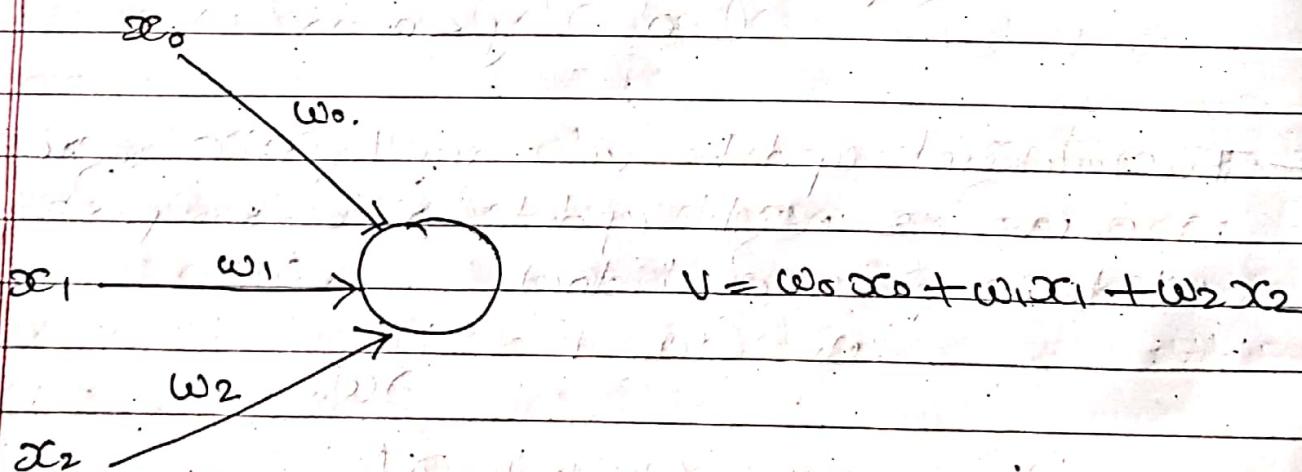
$$\frac{\partial L}{\partial w_0} = (\hat{y} - y) x^0 \quad (\text{where, } x^0 \text{ is always 1})$$

$$\frac{\partial L}{\partial w_1} = (\hat{y} - y) x^1$$

$$\frac{\partial L}{\partial w_2} = (\hat{y} - y) x^2$$

Thus, General weight update rule for perceptron learning is

$$w_i = w_i + \alpha (\hat{y} - y) x^i$$



* Gaussian Discriminant Analysis (GDA):

There are two broad Categories of Supervised Classification algorithms:

- Discriminative (Discriminant) Analysis Algorithm
- Generative Learning Algorithms

i) Discriminant Analysis Algorithm:

- It finds decision boundary between different classes.
- For e.g:- Logistic regression, perceptron learning, etc.
- For the test data that lies in one side of the boundary belongs to class 1 and the data point that lies in another side of decision boundary belongs to class 2.

ii) Generative Learning Algorithm:

- It finds the distribution of each class separately
- For e.g:- Naive Bayes classifier
- At the time of prediction, test data points is compared with each ^{data} distribution & belongs to the class with which it resembles the most.

* Naive Bayes classifier:

- It is based on Bayes theorem which is given as below:

$$P(H/x) = \frac{P(x/H) \cdot P(H)}{P(x)}$$

Posterior Probability \rightarrow Likelihood \rightarrow Prior Probability

- Bayes' theorem provides the way of calculating posterior probability from prior probabilities & likelihood.
 - Let, D be the dataset and c_1, c_2, \dots, c_m are m classes.
- now, Bayes' rule can be written as below:

$$p(c_i|x) = \frac{p(x|c_i) \cdot p(c_i)}{p(x)}$$

where, x is the input tuple containing n attributes i.e. $x = \{x_1, x_2, \dots, x_n\}$

The classifier predicts that tuple x belongs to class c_i having posterior probability conditional on x .

Since, $p(x)$ is same for all classes, we need to maximize $p(x|c_i) \cdot p(c_i)$

let x contains set of K attributes (i.e $x = \{x_1, x_2, x_3, \dots, x_K\}$) which are independent of each other.

now, the probability $p(x|c_i)$ can be written as below.

$$\begin{aligned} p(x|c_i) &= p(x_1|c_i) \times p(x_2|c_i) \dots \times p(x_K|c_i) \\ &= \prod_{j=1}^K p(x_j|c_i) \end{aligned}$$

$$\Rightarrow p(c_i|x) = p(c_i) \times \frac{\prod_{j=1}^K p(x_j|c_i)}{p(x)}$$

Finally, $p(x)$ is same for all classes so while making prediction, x is predicted to be in class C_i for which $P(C_i|x)$ is maximum.

Example:

Age	Income	Student	Credit_Rating	Buys_computer (class)
Youth	High	No	Fair	No
Youth	High	No	Excellent	No
Middle-Age	High	No	Fair	Yes
Senior	Medium	No	Fair	Yes
Senior	Low	Yes	Fair	Yes
Senior	Low	Yes	Excellent	No
Middle-Age	Low	Yes	Excellent	Yes
Youth	Medium	No	Fair	No
Youth	Low	Yes	Fair	Yes
Senior	Medium	Yes	Fair	Yes
Youth	Medium	Yes	Excellent	Yes
Middle-Age	Medium	No	Excellent	Yes
Middle-Age	High	Yes	Fair	Yes
Senior	Medium	No	Excellent	No

Consider the above training data & predict class labels if $X = S$. Age = Youth, Income = Medium, Student = Yes, Credit_Rating = Fair?

= Sol

$$P(\text{Buys} = \text{yes}) = 9/14 = 0.643 \quad P(\text{cr})$$

$$P(\text{Buys} = \text{No}) = 5/14 = 0.357$$

For Age : Youth

$$P(\text{Age} = \text{youth} | \text{Buys} = \text{yes}) = 2/9$$

$$P(\text{Age} = \text{youth} | \text{Buys} = \text{No}) = 3/5$$

For Income : Medium

$$P(\text{Income} = \text{medium} | \text{Buys} = \text{yes}) = 4/9$$

$$P(\text{Income} = \text{medium} | \text{Buys} = \text{No}) = 2/5$$

For student : Yes

$$P(\text{Student} = \text{Yes} | \text{Buys} = \text{yes}) = 6/9$$

$$P(\text{Student} = \text{Yes} | \text{Buys} = \text{No}) = 1/5$$

For credit Rating : Fair

$$P(\text{Credit Rating} = \text{Fair} | \text{Buys} = \text{yes}) = 6/9$$

$$P(\text{Credit Rating} = \text{Fair} | \text{Buys} = \text{No}) = 2/5$$

$$P(\text{Buys} = \text{yes} | x) = P(\text{Buys} = \text{yes}) \times \prod_{i=1}^n P(x_i | C_i)$$

$$= P(\text{Buys} = \text{yes}) \times P(\text{Age} = \text{youth} | \text{Buys} = \text{yes}) \times P(\text{Income} = \text{medium} | \text{Buys} = \text{yes}) \times P(\text{Student} = \text{Yes} | \text{Buys} = \text{yes}) \times P(\text{Credit Rating} = \text{Fair} | \text{Buys} = \text{yes})$$

$$= \frac{9}{14} \times \frac{2}{5} \times \frac{4}{9} \times \frac{6}{9} \times \frac{6}{9}$$

$$= 0.028.$$

Similarly,

$$p(\text{Buys} = \text{No} | x) = p(\text{Buys} = \text{No}) \times \prod_{j=1}^n p(x_j | C_j)$$

$$= p(\text{Buys} = \text{No}) \times p(\text{Age} = \text{youth} | \text{Buys} = \text{No}) \times \\ p(\text{Income} = \text{Mediocre} | \text{Buys} = \text{No}) \times p(\text{Student} = \text{yes} | \text{Buys} = \text{No}) \times p(\text{Credit Rating} = \text{Fair} | \text{Buys} = \text{No})$$

$$= \frac{5}{14} \times \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} \times \frac{2}{5}$$

$$= 0.007$$

Since, $p(\text{Buys} = \text{yes} | x) > p(\text{Buys} = \text{No} | x)$,

predicted class label for the given tuple is

"BuysComputer = yes"

* Variants of Naive Baye's classifier:

① Gaussian Naive Bayes:

- Gaussian Naive Bayes is used when features of data set have gaussian distribution.
- If features of dataset have continuous values, they follow gaussian distribution.
- In this case, probability of feature belonging to class C is estimated using gaussian probability distribution function as given below.

$$p(x_i | C) = \frac{1}{\sqrt{2\pi} \sigma_C} \exp \left(-\frac{(x_i - \mu_C)^2}{2\sigma_C^2} \right)$$

where, μ_C & σ_C are mean & standard deviation

of feature belonging to class C.

② Multinomial Naïve Bayes:

- It is used when input features follows multinomial distribution. It is primarily used in text classification.
- It calculates probability of features x_i belonging to class C as below:

$$p(x_i|c) = \frac{N_{ic} + \alpha}{N_c + \alpha n}$$

where,

N_{ic} is no. of occurrences of features x_i in class C.

N_c is the no. of samples belonging to class C.

α is smoothing factor (Normally $\alpha = 1$ is used).

n is no. of dimensions.

③ Bernoulli Naïve Bayes:

- It is used when input features follows Bernoulli distribution. When input features of dataset are binary value, they follow Bernoulli distribution.
- It calculates probability of feature x_i belonging to class C as below:

$$p(x_i|c) = p(x_i|c)x_i + (1 - p(x_i|c))(1 - x_i)$$

* Laplace Smoothing:

- It is the technique used to handle zero probability in Naive Bayes classifier.
- It estimates probability of features x_i belonging to class C as below: (Case of multinomial naive Bayes)

$$p(x_i | C) = \frac{N_{iC} + \alpha}{N_C + \alpha \cdot n}$$

where,

- N_{iC} is no. of occurrences of feature x_i in class C.
- N_C is the no. of samples belonging to class C.
- α is smoothing factor (Normally $\alpha = 1$ is used)
- n is no. of dimensions

* Example:

Consider the following data set:

Email	Class
Send us your password	Spam
Send your Account	Spam
Review your Account	Ham
Review your password	Ham
Review our profile	Spam
Send us money	Spam
Review your profile	?

= Sol

Assume $\alpha = 1$

Vocabulary = {send, us, your, password, account, review, our, profile, money?}

$$P(\text{spam}) = 4/6$$

$$P(\text{Ham}) = 2/6$$

$$\therefore P(x|c) = \frac{N_c + 1}{N_c + N_H}$$

① Send:

$$P(\text{Send} | \text{spam}) = \frac{3+1}{4+1 \times 9} = \frac{4}{13}$$

$$P(\text{Send} | \text{Ham}) = \frac{0+1}{2+1 \times 9} = \frac{1}{11}$$

② US:

$$P(\text{US} | \text{spam}) = \frac{2+1}{4+1 \times 9} = \frac{3}{13}$$

$$P(\text{US} | \text{Ham}) = \frac{0+1}{2+1 \times 9} = \frac{1}{11}$$

③ Yoces:

$$P(\text{Yoces} | \text{spam}) = \frac{2+1}{4+1 \times 9} = \frac{3}{13}$$

$$P(\text{Yoces} | \text{Ham}) = \frac{2+1}{2+1 \times 9} = \frac{3}{11}$$

④ password:

$$P(\text{password} | \text{spam}) = \frac{1+1}{4+1 \times 9} = \frac{2}{13}$$

$$P(\text{password} | \text{Ham}) = \frac{1+1}{2+1 \times 9} = \frac{2}{11}$$

⑤ Account :

$$P(\text{Account} | \text{Spam}) = \frac{1+1}{4+1 \times 9} = \frac{2}{13}$$

$$P(\text{Not Account} | \text{Ham}) = \frac{1+1}{2+1 \times 9} = \frac{2}{11}$$

⑥ review :

$$P(\text{Review} | \text{Spam}) = \frac{1+1}{4+1 \times 9} = \frac{2}{13}$$

$$P(\text{Review} | \text{Ham}) = \frac{2+1}{2+1 \times 9} = \frac{3}{11}$$

⑦ user :

$$P(\text{User} | \text{Spam}) = \frac{1+1}{4+1 \times 9} = \frac{2}{13}$$

$$P(\text{User} | \text{Ham}) = \frac{0+1}{2+1 \times 9} = \frac{1}{11}$$

⑧ profile :

$$P(\text{Profile} | \text{Spam}) = \frac{1+1}{4+1 \times 9} = \frac{2}{13}$$

$$P(\text{Profile} | \text{Ham}) = \frac{0+1}{2+1 \times 9} = \frac{1}{11}$$

⑨ money :

$$P(\text{Money} | \text{Spam}) = \frac{1+1}{4+1 \times 9} = \frac{2}{13}$$

$$P(\text{Money} | \text{Ham}) = \frac{0+1}{2+1 \times 9} = \frac{1}{11}$$

now,

$$\begin{aligned}
 p(\text{spam} | \text{"Review your profile"}) &= p(\text{spam} | [0, 0, 1, 0, 0, 1, 0, 1]) \\
 &= p([0, 0, 1, 0, 0, 1, 0, 1] | \text{spam}) \times p(\text{spam}) \\
 &= \left(1 - \frac{4}{13}\right) \times \left(1 - \frac{3}{13}\right) \times \left(\frac{3}{13}\right) \times \left(1 - \frac{2}{13}\right) \times \left(1 - \frac{2}{13}\right) \times \frac{2}{13} \\
 &\quad \times \left(1 - \frac{2}{13}\right) \times \frac{2}{13} \times \left(1 - \frac{2}{13}\right) \times \frac{4}{6} \\
 &= 0.000994
 \end{aligned}$$

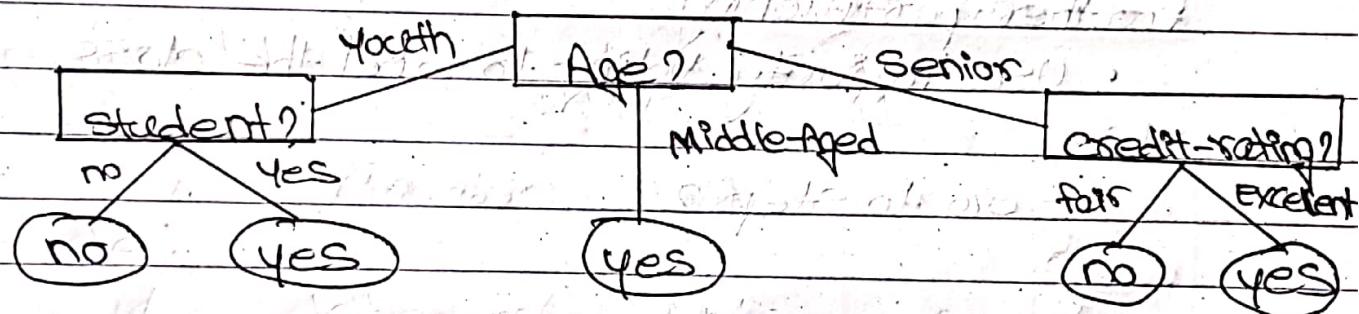
$$\begin{aligned}
 p(\text{Ham} | \text{"Review your profile"}) &= p(\text{Ham} | [0, 0, 1, 0, 0, 1, 0, 1]) \\
 &= p([0, 0, 1, 0, 0, 1, 0, 1] | \text{Ham}) \times p(\text{Ham}) \\
 &= \left(1 - \frac{1}{11}\right) \times \left(1 - \frac{1}{11}\right) \times \left(\frac{3}{11}\right) \times \left(1 - \frac{2}{11}\right) \times \left(1 - \frac{2}{11}\right) \times \frac{3}{11} \times \\
 &\quad \left(1 - \frac{1}{11}\right) \times \frac{1}{11} \times \left(1 - \frac{1}{11}\right) \times \frac{2}{6} \\
 &= 0.00103
 \end{aligned}$$

since, $P(\text{Ham} | \text{"Review your profile"}) > P(\text{spam} | \text{"Review your profile"})$.

predicted class label for the given email message is "Ham".

* Decision Tree classifier:

- It is the supervised classification algorithm that constructs the decision tree from the given training dataset.
- Decision tree is the tree structure where internal nodes (non-leaf node) denotes a test on an attribute, branches or edges represents outcomes of tests, and leaf nodes (terminal nodes) holds class labels.
- In order to make prediction for a tuple, the attributes of a tuple are tested against the decision tree. A path is traced from the root to a leaf node which determines the predicted class for that tuple.



A decision tree indicating whether a customer is likely to purchase a computer.

- i) Class-label Yes : The customer likely to buy a computer
- ii) Class-label No : The customer is unlikely to buy a computer.

* Algorithm for Constructing Decision Tree classifier:

Constructing a decision tree uses greedy algorithm. Tree is constructed in a top-down divide-and-conquer manner.

1. Initially, all the training tuples are considered in root nodes.
2. Tuples are partitioned recursively based on selected attributes.
3. If all samples of the given node belongs to the same class.
 - Label the class.
4. Else if there are no remaining attributes for further partitioning.
 - Use majority voting to label the class.
5. ELSE
 - Go to step 2
 - These are many variations of decision tree algorithms. Some of them are: ID 3 (Iterative Dichotomiser 3), C4.5 (successor of ID 3), CART (Classification and Regression Tree) etc.
 - There are different attribute selection measures used by decision tree classifier. Some of them are: Information Gain, Gain Ratio, Gini Index etc.

* ID3 Decision Tree:

- ID3 stands for iterative dichotomiser 3. It uses top-down greedy approach to build decision tree model.
- This algorithm uses information gain as attribute selection measure and selects the attributes with the highest information gain.
- Information gain measures reduction in entropy after data transformation.
- Entropy is the measure of homogeneity of the sample. Entropy or expected information of dataset D is calculated by using equation below:

$$E(D) = - \sum_{i=1}^m P_i \log_2 P_i$$

where, P_i is the probability of a sample belonging to class C_i and

m is the number of classes.

$$P_i = \frac{|C_i, D|}{|D|}$$

where, $|C_i, D|$ is the total number of samples in D belonging to class C_i and $|D|$ is the total number of samples in training dataset.

- Suppose dataset D is partitioned on the basis of attribute A having V distinct values. Thus, dataset D is partitioned into V partitions $\{D_1, D_2, \dots, D_V\}$

hence, the total entropy of the data set after partitioning on A can be calculated as below:

$$EA(D) = \sum_{i=1}^{|D|} |D_i| E(D_i)$$

- Information gain is reduction in entropy after partitioning - Thus, information after partitioning on A can be calculated as below.

$$\text{Gain}(A) = E(D) - EA(D)$$

- ID3 selects attributes with highest information for partitioning.

* Example:

Age	Income	Student	Credit Rating	Buys Computer
Youth	Low	Yes	Good	Yes
Youth	Low	No	Good	No
Middle Age	Medium	Yes	Good	Yes
Senior	High	Yes	Bad	No
Senior	High	No	Bad	No
Middle Age	Medium	Yes	Bad	No
Youth	Low	Yes	Bad	No
Youth	Low	No	Bad	No
Senior	High	Yes	Bad	No
Youth	Low	Yes	Good	Yes
Middle Age	Medium	Yes	Good	Yes
Senior	High	No	Good	No

Age	Income	Student	Credit Rating	(Class) Buys Computer
Youth	High	No	Fair	No
Youth	High	No	Excellent	No
Middle-Age	High	No	Fair	Yes
Senior	Medium	No	Fair	Yes
Senior	Low	Yes	Fair	Yes
Senior	Low	Yes	Excellent	No
Middle-Age	Low	Yes	Excellent	Yes
Youth	Medium	No	Fair	No
Youth	Low	Yes	Fair	Yes
Senior	Medium	Yes	Fair	Yes
Youth	Medium	Yes	Excellent	Yes
Middle-Age	Medium	No	Excellent	Yes
Middle-Age	High	Yes	Fair	Yes
Senior	Medium	No	Excellent	No

- Construct decision tree from above data using information gain.
- predict class label of the tuple : $x = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit-rating} = \text{fair})$

= S0 Entropy of the data set before partitioning is given as below :

We have 9 positive & 5 negative class out of 14 class
So, $S(9+, 5-) :$

$$E(D) = - \sum_{i=1}^m P_i \log_2 P_i$$

$$= - \frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.41 + 0.53$$

$$= 0.94$$

① Attribute = Age

$$\text{Syouth } [2+, 3-] \quad \text{Entropy (Syouth)} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.971$$

$$\text{Smiddle age } [4+, 0-] \quad \text{Entropy (Smiddle age)} = -\frac{4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7}$$

$$= 0$$

[In case of either
+ or - = 0 or - = 0 Ans]

$$\text{Ssenior } [3+, 2-] \quad \text{Entropy (Ssenior)} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.971$$

$$\text{Gain } (S, \text{Age}) = E(S) - \sum_{i=1}^{|D|} |D_i| \cdot E(D_i)$$

$$= \text{Entropy}(S) - \left(\frac{5}{14} \text{Entropy}(\text{Syouth}) + \frac{4}{14} \text{Entropy} \right.$$

$$\left. (\text{Smiddle age}) + \frac{5}{14} \text{Entropy}(\text{Ssenior}) \right)$$

$$\text{Gain } (S, \text{Age}) = 0.94 - \left(\frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 \right)$$

$$= 0.94 - 0.694$$

$$= 0.246$$

② Attribute = "Income"

$$S_{\text{High}} [2+, 2-] \text{ Entropy} (S_{\text{High}}) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4}$$

$$= 1.0 \quad [\because \text{in case of equal + & - entropy is 1}]$$

$$S_{\text{Medium}} [4+, 2-] \text{ Entropy} (S_{\text{medium}}) = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6}$$

$$= 0.9183$$

$$S_{\text{Low}} [3+, 1-] \text{ Entropy} (S_{\text{low}}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4}$$

$$\text{Gain}(S, \text{Income}) = E(S) - \sum_{i=1}^3 P(D_i) E(D_i)$$

$$= E(S) - \left(\frac{4}{14} \times 1 + \frac{6}{14} \times 0.9183 + \frac{4}{14} \times 0.811 \right)$$

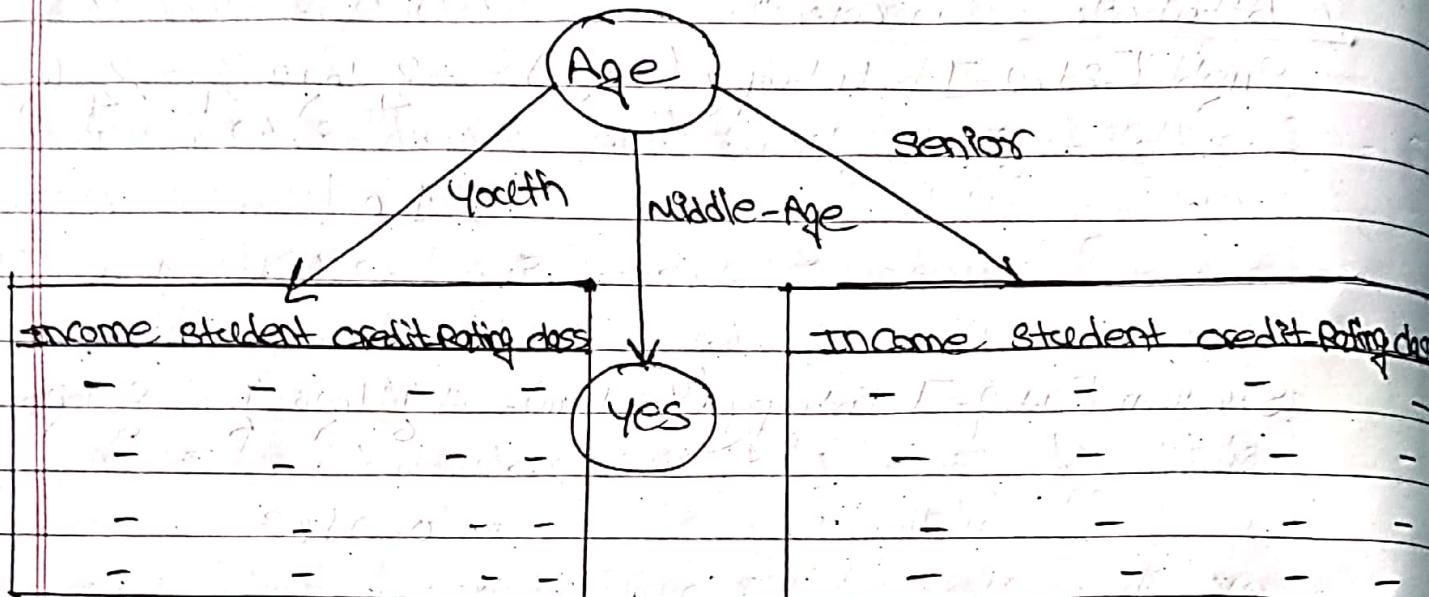
$$= 0.94 - 0.904$$

Similarly,

$$\text{Gain}(\text{student}) = 0.151$$

$$\text{Gain}(\text{credit rating}) = 0.048$$

Because age has the highest information gain among the attributes, it is selected as the splitting attribute. The decision tree now looks like below:



- Notice that the tuples falling into the partition for age > middle-age all belongs to the same class. Because they all belongs to class "yes" a leaf should therefore be created at the end of this branch & labeled "Yes".

now, consider rows where, only youth is represent we get: (case of left children)

Age	Income	Student	Credit Rating	Class
Youth	High	No	Fair	No
Youth	High	No	Excellent	No
Youth	Medium	No	Fair	No
Youth	Low	Yes	Fair	Yes
Youth	Medium	Yes	Excellent	Yes

Repeat the same process as above we get,

Entropy of the dataset before partitioning is given as

$$S(2+, 3-) \quad E(D) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.528 + 0.442$$

$$= 0.97$$

Attribute = "Income"

$$S(0+, 2-) \quad E(S_{High}) = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2}$$

$$= 0 \quad [\because \text{In case of any one is } 0]$$

$$S_{Medium}(1+, 1-) \quad E(S_{Medium}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 1$$

$$S_{Low}(1+, 0-) \quad E(S_{Low}) = -\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1}$$

$$= 0$$

now,

$$\text{Gain}(S_{yacht}, \text{Income}) = \text{Entropy}(D) - \sum_{i=1}^{|D|} |D_i| E(D_i)$$

$$= E(D) - \left(\frac{2}{5} \times 0 + \frac{2}{5} \times 1 + \frac{1}{5} \times 0 \right)$$

$$= 0.97 - 0.4$$

$$= 0.57$$

✓

Attribute = "Student"

$$\textcircled{1} S_{No}[0+, 8-] \quad E(S_{No}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0$$

$$\textcircled{2} S_{yes}[2+, 0-] \quad E(S_{yes}) = 0$$

So,

$$\begin{aligned} \text{Gain}(S_{youth}, \text{student}) &= E(D) - \sum_{i=1}^V p_i E(D_i) \\ &= 0.97 \left(\frac{3}{5} \times 0 + \frac{2}{5} \times 0 \right) \\ &= 0.97 \end{aligned}$$

Attribute = "credit rating"

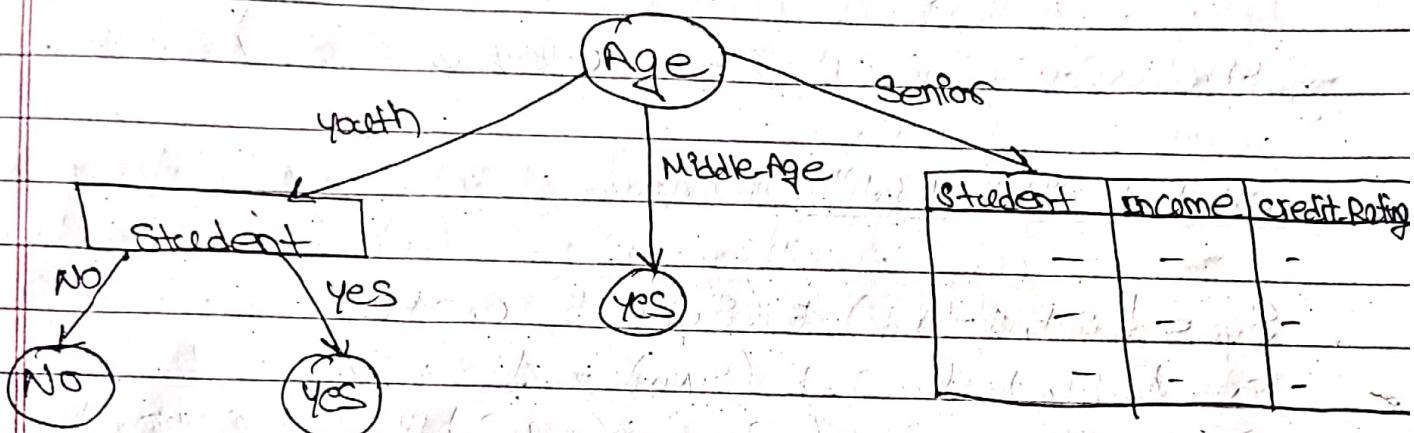
$$\begin{aligned} \textcircled{1} S_{Fair}[2+, 1-] \quad E(S_{Fair}) &= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ &= 0.9183 \end{aligned}$$

$$\textcircled{2} S_{Excellent}[1+, 1-] \quad E(S_{Excellent}) = 1 \quad [\because \text{both are same}]$$

So,

$$\begin{aligned} \text{Gain}(S_{youth}, \text{credit Rating}) &= E(D) - \sum_{i=1}^V p_i E(D_i) \\ &= 0.97 - \left(\frac{3}{5} \times 0.9183 + \frac{2}{5} \times 1 \right) \\ &= 0.0192 \end{aligned}$$

now, Attribute "Student" has the highest information gain among the other attributes we use as splitting attributes. So, the decision tree we get is.



now, Some method is repeated for the ~~left~~ Right children i.e. "Senior". So, Consider only rows whose Senior appears i.e.

Age	Income	Student	Credit Rating	Student	Credit Rating	Class
Senior	Medium	Fair	Fair	No	?	yes
Senior	Low	Fair	Fair	Yes	?	yes
Senior	Low	Excellent	Excellent	Yes	?	No
Senior	Medium	Fair	Fair	Yes	?	yes
Senior	Medium	Excellent	Excellent	No	?	No

now,

Entropy of dataset before partitioning is given as:

$$S(3+, 2-) \quad E(S) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$E(\text{Senior}) = 0.97$$

Attribute "Income"

S_{Medium} (2+, 1-)

$$E(\text{Smedium}) = 0.9183$$

S_{Low} (1+, 1-)

$$E(\text{Slow}) = 1$$

$$\text{Gain}(\text{Senior}, \text{Income}) = G_1(D) - \sum_{i=1}^3 \frac{|D_i|}{|D|} E(D_i)$$

$$= 0.97 - \left(\frac{3}{5} \times 0.9183 + \frac{2}{5} \times 1 \right)$$

$$= 0.0192 \checkmark$$

Attribute "Student"

$$S_{\text{Yes}} [2+, 1-] \quad E(S_{\text{Yes}}) = 0.9183$$

$$S_{\text{No}} [1+, 1-] \quad E(S_{\text{No}}) = 1$$

So,

$$\text{Gain}(\text{Senior}, \text{Student}) = 0.0192 \checkmark$$

Attribute "Credit Rating"

$$S_{\text{Fair}} [3+, 0-] \quad E(S_{\text{Fair}}) = 0$$

$$S_{\text{Excellent}} [0+, 2-] \quad E(S_{\text{Excellent}}) = 0$$

So,

$$\text{Gain} = 0.97 \checkmark$$

now, Attribute "credit rating" has high information Gain value, we use it as a splitting attribute so we get the decision tree as below.

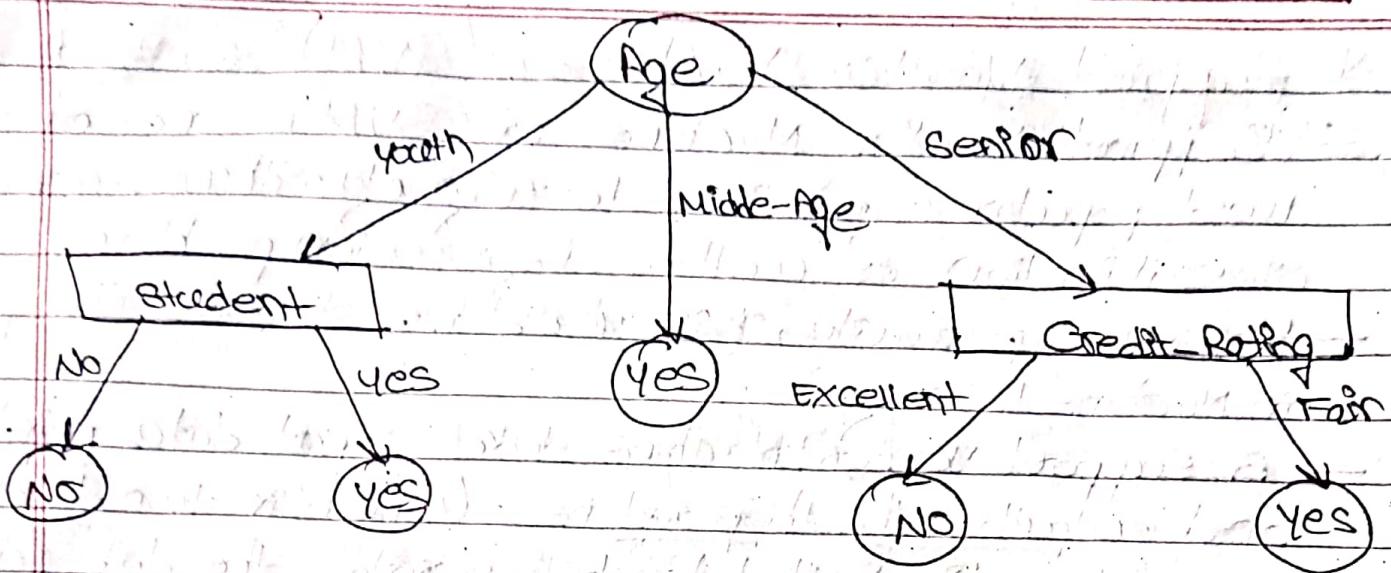
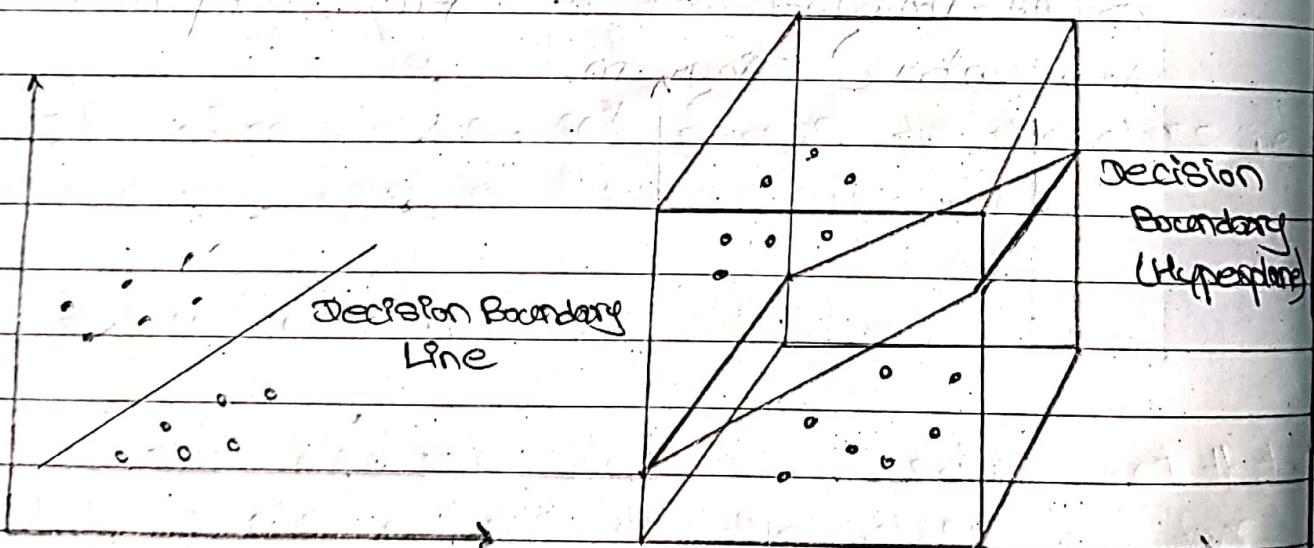


Fig:- Final decision tree using ID3 Algo.

Predicted class level of the people
 $X = (\text{Age} = \text{youth}, \text{income} = \text{medium}, \text{Student} = \text{yes}, \text{credit rating} = \text{fair})$ is $\text{Burg Computer} = \text{yes}$
 from above diagram.

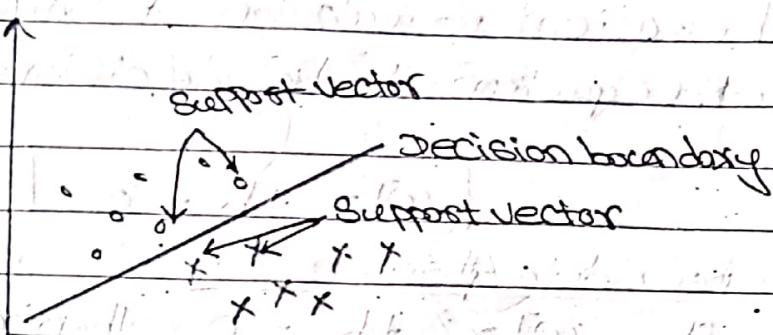
* Support Vector Machine (SVM):

- Support Vector Machine or SVM is one of the most popular supervised learning algorithms, used for classification as well as Regression problems.
- However, primarily, it is used for classification problems in Machine Learning.
- A Support Vector Machine takes input data points and accepts the hyperplane (which in two dimensions it's simply a line) that best separates the data points into two classes.
- This line or hyperplane is the decision boundary. Any data points that falls to one side of it is classified in one class and, the data points that falls to the other of it is classified in another class.



- Compared to newer algorithms like neural networks, they have two main advantages: higher speed and better performance with a limited number of samples (in the thousands).

- Support vectors are data points that are closest to the hyperplane and influence the position and orientation of the hyperplane.



- To separate the two classes of data points, there are many possible hyperplanes that could be chosen. SVM algorithm selects the optimal hyperplane by choosing hyperplane with largest margin. Such hyperplane is called Maximum marginal hyperplane (MMH).

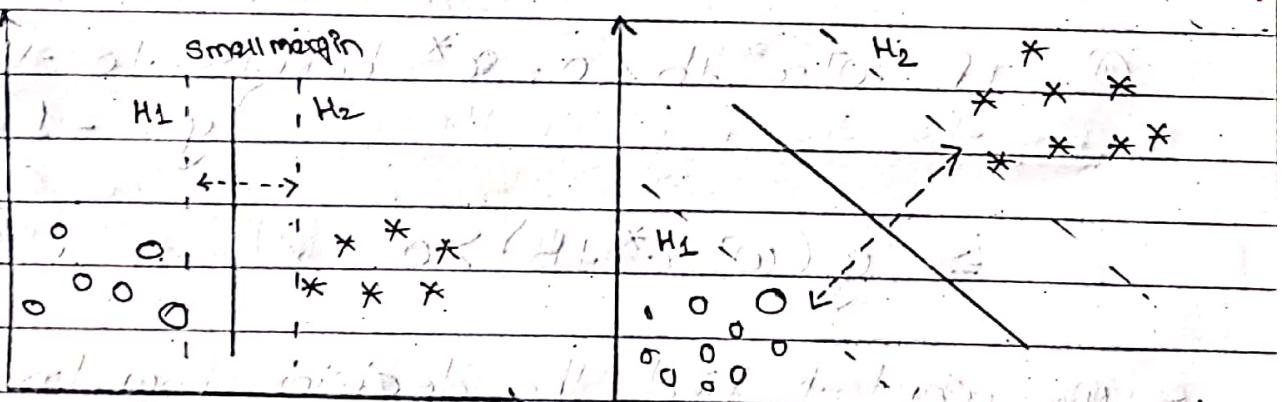


Fig.: Decision boundary with small margin

Fig.: Decision boundary with large Margin

- Let H_1 & H_2 be two planes that pass through support vectors and parallel to the hyperplane of decision boundary. Margin is the distance between H_1 & H_2 such that distance between H_1 & decision boundary

must be equal to the distance between H₂ & the decision boundary.

* Mathematical formulation of SVM:

Let, equation of the decision boundary is

$$\omega^T x + b = 0 \quad \dots \textcircled{1}$$

for new data point x^*

If $\omega^T x^* + b = 0$, then x^* is on the decision boundary

If $\omega^T x^* + b > 0$, then x^* is above the decision boundary

If $\omega^T x^* + b < 0$, then x^* is below the decision boundary

For classification,

① If $\omega^T x^* + b > 0$; x^* belongs to the class +ve i.e. $y = +1$

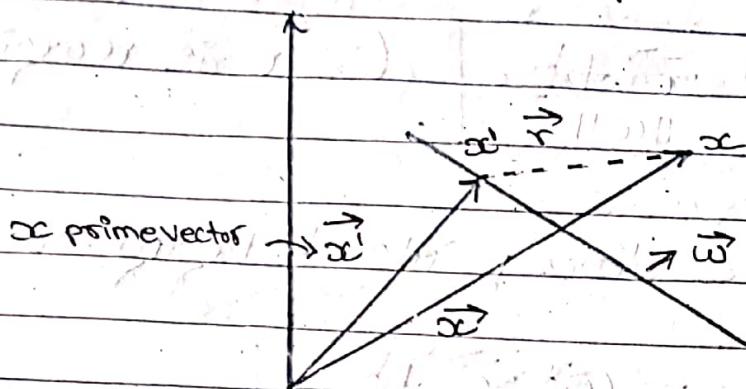
② If $\omega^T x^* + b < 0$; x^* belongs to -ve class i.e. $y = -1$

$$\Rightarrow y(\omega^T x^* + b) > 0$$

- We want to find the decision boundary such that all data points belongs to the class +ve for from the decision boundary as well as the points belonging to -ve class are also far from the decision boundary so that class label for every data point can be predicted with higher confidence i.e. we want $|\omega^T x^* + b|$ is large.

\Rightarrow We want $y(\omega^T \vec{x}^* + b)$ is large.

now, for the point belonging to the class,



$$\vec{r} = \vec{x} - \vec{x}' \quad \text{--- (6)}$$

for, x' , $\omega^T \vec{x}' + b = 0$ --- (3) x' lies on decision boundary.

Also,

\vec{r} is always parallel to \vec{w}

$$\text{let } |\vec{r}| = r$$

$$\Rightarrow \vec{r} = r \cdot \vec{w} = r \cdot \frac{\vec{w}}{\|\vec{w}\|} \quad [\because \vec{w} \text{ is unit vector}]$$

from eqn (6) we have

$$\vec{x}' = \vec{x} - \vec{r} = \vec{x} - \frac{r \cdot \vec{w}}{\|\vec{w}\|}$$

Substituting this value in eqn (3), we get

$$\omega^T \left(\vec{x} - \frac{r \cdot \vec{w}}{\|\vec{w}\|} \right) + b = 0$$

$$\Rightarrow \omega^T \vec{x} - \frac{r \cdot \vec{w} \cdot \omega^T}{\|\vec{w}\|} + b = 0$$

since, $\omega^T \vec{w} = \|\omega\|^2$

$$\Rightarrow \omega^T \vec{x} - r\|\omega\| + b = 0$$

$$\Rightarrow r = \frac{\omega^T \vec{x} + b}{\|\omega\|} \quad (\because r \text{ is margin})$$

Similarly, for the data points belonging to -ve class

$$r = \frac{-\omega^T \vec{x} + b}{\|\omega\|}$$

Thus, in general margin is given as below,

$$r = \frac{y(\omega^T \vec{x} + b)}{\|\omega\|}$$

① For the +ve point on H₁,

$$r = \frac{\omega^T \vec{x} + b}{\|\omega\|} = \frac{1}{\|\omega\|} \text{ or } \frac{C}{\|\omega\|}$$

② For the -ve point on H₂,

$$r = \frac{-\omega^T \vec{x} + b}{\|\omega\|} = \frac{-1}{\|\omega\|} \text{ or } \frac{-C}{\|\omega\|}$$

$$\Rightarrow \text{Total margin} = \frac{C}{\|\omega\|} + \frac{-C}{\|\omega\|} = \frac{2C}{\|\omega\|} = 2 \text{ when } C=1$$

Our goal is to maximize above margin. Maximizing $\frac{1}{2} \|w\|^2$ is equivalent to minimizing $\|w\|$. Thus, we can formulate following optimization problem as:

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

$$\text{Subject to } y_i(w^T x_i + b) \geq 1, \quad i=1,2,3,\dots,n$$

Above optimization problem can be solved using quadratic programming.

* Functional and Geometric Margins of SVM:

We can write SVM classifier as :

$$h_{w,b}(x) = g(w^T x + b)$$

where,

$$g(z) = 1 \quad \text{if } z > 0$$

$$g(z) = -1 \quad \text{otherwise}$$

Given a training example (x_i, y_i) , functional margins of the decision boundary parameterized by (w, b) is given as below:

$$x_i = y_i (w^T x_i + b)$$

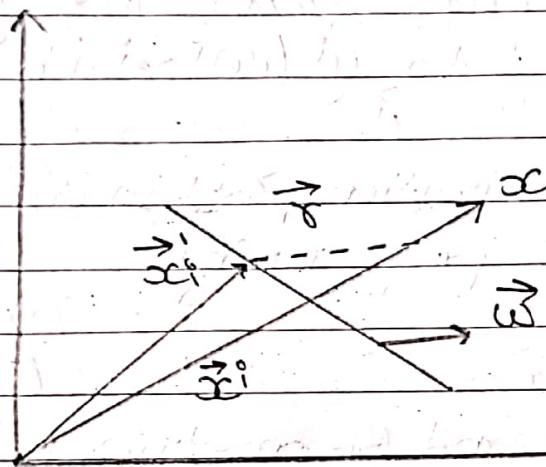
If we replace w by $2w$ & b by $2b$ then

$$g(w^T x + b) = g(2w^T x + 2b)$$

- We can make functional margin arbitrarily large

without changing meaning. This is the drawback of functional margin.

To solve this problem geometric margin is introduced for a the point \vec{x}_i^* .



$$\vec{r} = \vec{x}_i - \vec{x}_i'$$

$$\Rightarrow \vec{x}_i' = \vec{x}_i + \vec{r} = \vec{x}_i + r_i \vec{w} \quad \because \frac{\vec{w}}{\|\vec{w}\|}$$

where $|\vec{r}| = r_i$

since, \vec{x}_i' is point on decision boundary,

$$\Rightarrow \vec{w}^T \vec{x}_i' + b = 0$$

$$\Rightarrow \vec{w}^T (\vec{x}_i + r_i \vec{w})$$

$$\Rightarrow \frac{\vec{w}^T (\vec{x}_i + r_i \vec{w})}{\|\vec{w}\|} + b = 0$$

$$\Rightarrow r_i = \frac{w^T x_i + b}{\|w\|} = \left(\frac{w}{\|w\|} \right)^T x_i + \frac{b}{\|w\|}$$

i.e. $r_i = \left(\frac{w}{\|w\|} \right)^T x_i + \frac{b}{\|w\|}$

Thus, geometric margin of the decision boundary parameterized by (w, b) is given as below:

$$r_i = y_i \left(\left(\frac{w}{\|w\|} \right)^T x_i + \frac{b}{\|w\|} \right)$$

* Lagrange Duality:

- In mathematical optimization theory, duality is the principle that optimization problem can be viewed from either of two perspective that is the primal problem and the dual problem.
- If the primal problem is minimization problem then the dual problem is maximization problem & vice versa.
- Lagrangian duality theory refers to the way to solve an primal optimization problem by looking at dual optimization problem.

* Karush - Kuhn - Tucker (KKT) conditions:

- KKT conditions are tests on first derivatives and other Regularity Conditions that needs to be satisfied for solving a non-linear optimization problem.
- KKT conditions for solving SVM optimization problem are:

$$\frac{\partial L}{\partial w_i} = 0$$

$$\frac{\partial L}{\partial b_i} = 0$$

$$\begin{cases} \sum_i c_i g_i(w) = 0 \\ g_i(w) \leq 0 \\ c_i \geq 0 \end{cases}$$

where,

L is Lagrangian duality function.

* Solving SVM optimization problem:

SVM optimization problem is given by

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

$$\text{s.t.c } y_i(w^T x_i + b) \geq 1$$

for $i = 1, 2, \dots, m$

Above, optimization problem can also be written as below:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. c } q_i(w) = -y_i(w^T x_i + b) + 1 \leq 0$$

Lagrangian function for the given optimization problem is

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) + 1]$$

①

where, α is Lagrange multiplier.
now,

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^m \alpha_i y_i x_i = 0$$

$$\Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m \alpha_i y_i = 0$$

Substituting values of w in eqn ①, we get:

$$\frac{1}{2} \sum_{i=1}^m \alpha_i y_i x_i \cdot \sum_{j=1}^m \alpha_j y_j x_j - \sum_{i=1}^m \alpha_i [y_i (\sum_{j=1}^m \alpha_j y_j x_j + b) + 1]$$

Solving this we get,

$$L(w, b, \alpha) = \sum \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j x_i x_j$$

Thus, the task is to solve the following optimization:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j x_i x_j$$

$$\text{s.t. } \alpha_i \geq 0$$

$$\alpha_i y_i = 0$$

multiplier relaxation of the vector α is given by $\alpha^* = \alpha + \lambda y$

where λ is the Lagrange multiplier for the constraint $\alpha_i y_i = 0$.

GP [13] - If all $\alpha_i > 0$, then $\lambda = 0$

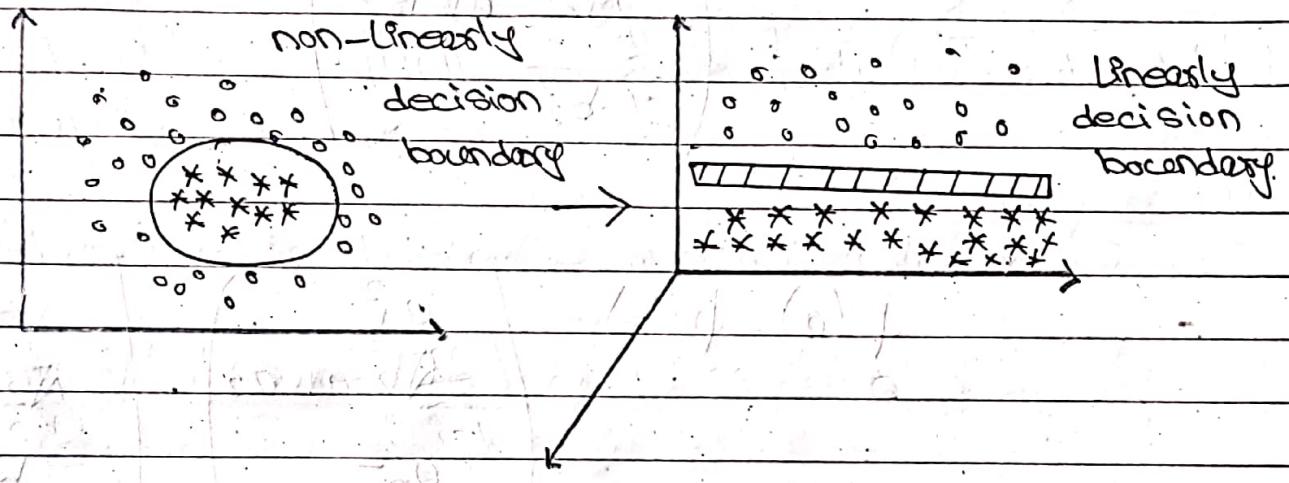
which means $\alpha^* = \alpha$ and $\alpha^* \neq \alpha$

so $\alpha^* \neq \alpha$ which contradicts the optimality of α

thus contradiction implies $\alpha_i = 0$ for all i

* SVM Kernel :

- SVM uses a set of mathematical functions which are called Kernels or Kernel functions.
- These functions are used to transform non-linearly separable low-dimensional data points to linearly separable high dimensional data points i.e.



Definition: Kernel function is a mathematical function that takes vector in original space as input & returns dot product of vector in transformed (High-dimensional) Space.
i.e.

$$K(a, b) = (a \cdot b)^2$$

$$= \phi(a) \cdot \phi(b)$$

Consider, the following Kernel function.

$$K(a, b) = \phi(a) \cdot \phi(b)$$

where,

$$\phi(x_1) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 \cdot x_2 \\ x_2^2 \end{pmatrix}$$

here,

$$\phi(a) = \begin{pmatrix} a_1^2 \\ \sqrt{2} a_1 a_2 \\ a_2^2 \end{pmatrix}$$

$$\phi(b) = \begin{pmatrix} b_1^2 \\ \sqrt{2} b_1 \cdot b_2 \\ b_2^2 \end{pmatrix}$$

now,

$$\phi(a) \cdot \phi(b) = \begin{pmatrix} a_1^2 \\ \sqrt{2} a_1 \cdot a_2 \\ a_2^2 \end{pmatrix} \cdot \begin{pmatrix} b_1^2 \\ \sqrt{2} b_1 \cdot b_2 \\ b_2^2 \end{pmatrix}$$

$$= a_1^2 \cdot b_1^2 + 2 a_1 \cdot b_1 \cdot a_2 \cdot b_2 + b_2^2 \cdot a_2^2$$

$$= (a_1 b_1 + a_2 b_2)^2$$

$$(a \cdot b)^2 = \left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \cdot \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \right)^2$$

$$= (a \cdot b)^2$$

* Example: Let say, we have two data points.

$$a = (1, 2), b = (2, 3)$$

now,

$$\phi a = (1, \sqrt{2} \times 2, 4)$$

$$\phi b = (4, \sqrt{2} \times 6, 9)$$

$$\phi a \cdot \phi b = (1, \sqrt{2} \times 2, 4) \cdot (4, \sqrt{2} \times 6, 9)$$

$$\begin{aligned} \text{LHS} &= (4 + 12 \times 2 + 36) \\ &= 64 \end{aligned}$$

using Kernel function, we have

$$k(a, b) = (a \cdot b)^2$$

$$\text{or constant value} = ((1, 2) \cdot (2, 3))^2$$

$$= (2 + 6)^2$$

$$= 64$$

Conclusion: Computing dot product of vectors in transformed space is computationally expensive. Kernel functions allows us to do this in computationally efficient manner.

* Types of Kernel:

① Linear Kernel:

$$K(a, b) = (a \cdot b)^2$$

② Polynomial Kernel:

$$K(a, b) = (a \cdot b + c)^d$$

(hyperparameter)

where, d is degree of polynomial.

③ Radial basis function kernel:

$$K(a, b) = \exp(-\gamma \|a - b\|^2)$$

R is parameter whose values ranges between 0 to 1.

④ Gaussian kernel:

$$K(a, b) = \exp\left(-\frac{\|a - b\|^2}{2\sigma^2}\right)$$

where, σ is bandwidth.

⑤ Sigmoid kernel:

It is equivalent to two layer perceptron model of neural network & is given as below:

$K(a,b) = \tanh(\alpha a^T b + c)$ where, α lies between 0 & 1.

⑥ Exponential Kernel:

$$K(a,b) = \exp\left(-\frac{\|a-b\|}{2\sigma^2}\right)$$

* Kernel Mercer's Theorem:

- Kernel matrix is symmetric matrix that is obtained by applying kernel function K in every pair of points in the data set.

i.e., $K_{ij} = K(x_i, x_j)$ where, $i = 1, 2, \dots, n$

$$K_{ij} = K(x_i, x_j) \quad i=1, 2, \dots, n \quad j=1, 2, \dots, n$$

$$\text{or, } K_{ij} = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2), & \dots & K(x_1, x_n) \\ K(x_2, x_1), & K(x_2, x_2), & \dots & K(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ K(x_n, x_1), & K(x_n, x_2), & \dots & K(x_n, x_n) \end{bmatrix}$$

- Mercer's theorem determines which functions can be used as a kernel functions.
- In mathematics, Mercer's theorem states that necessary & sufficient condition for a kernel K

to be valid is that it's kernel matrix is symmetric positive semi-definite.

* The Sequential Minimization optimization (SMO) Algorithm:

- SMO algorithm can be used to solve SVM dual problem efficiently.

* Coordinate Ascent Algorithm:

- It is method of solving unconstrained optimization problem of the form:

$$\max_w w(\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_n)$$

Algorithm:

Loop until Convergence

for $i = 1$ to n

$$\alpha_i = \arg \max_{\alpha} w(\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_{i-1}, \alpha_i, \alpha_{i+1}, \dots, \alpha_n)$$

* SMO Algorithm:

The SVM dual problem we have to solve is

$$w(\alpha) = \underset{\alpha}{\text{max}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (x_i \cdot x_j) \quad \text{--- (1)}$$

Such that, $\alpha_i > 0$ — (2)

$$\sum_{i=1}^n \alpha_i y_i = 0 - (3)$$

Coordinate ascent algorithm discussed above cannot be used to solve SVM dual problem because if we optimize α_1 by keeping α_2 to α_n fixed, it leads to

$$\alpha_1 y^1 = - \sum_{i=2}^n \alpha_i y^i$$

$$\Rightarrow \alpha_1 = \frac{1}{n-1} \sum_{i=2}^n \alpha_i y^i \text{ but } \alpha_i > 0 \quad \text{and } y^1 \neq 0$$

This violates constraints (2)

Thus, we must update at least two α_i simultaneously in order to keep constraint satisfying.

Now, SMO algorithm for solving SVM dual problem is given as below:

Step 1: Repeat until convergence

Step 2: Select pair of α_i and α_j

Step 3: Optimize $w(\alpha)$ w.r.t. α_i & α_j by keeping α_k fixed, $k = 1, 2, 3, \dots, n$, $k \neq i, j$

Step 4: Update α_i & α_j by checking for溢出 (overflow) or underflow (underflow)

Step 5: Check for convergence criterion, if satisfied, go to step 1, otherwise go to step 2

Step 6: Stop, repeat the process until all constraints are satisfied