

Name: Dhaval Kumar Vijay Kumar Patel

classmate

Date

Page

M.Sc CS

TA-2 Assignment

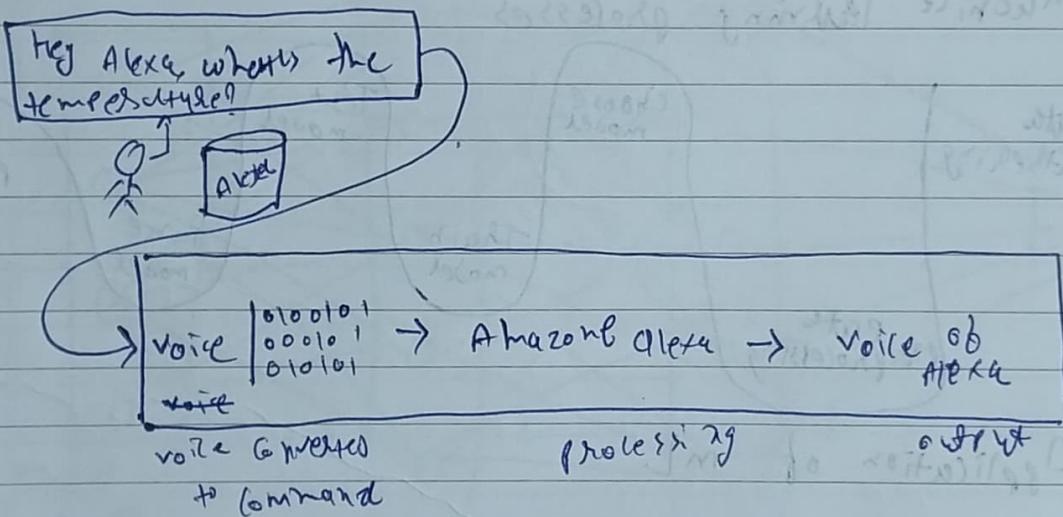
Date of submission
02/11/2022

Artificial Intelligence

a.) Illustrate with an example the difference between Machine Learning, Deep Learning and Artificial Intelligence.

→ b) Artificial Intelligence.

Artificial Intelligence (AI) is the process of imparting data, information and human intelligence to machine. The main goal of AI is to develop self-reliant machine that can think and act like human. These machines can mimic human behavior and perform tasks by learning and problem-solving. Most of the AI systems simulate natural intelligence to solve complex problems. Basically Artificial Intelligence measure Ability of a machine to imitate intelligent human behavior.



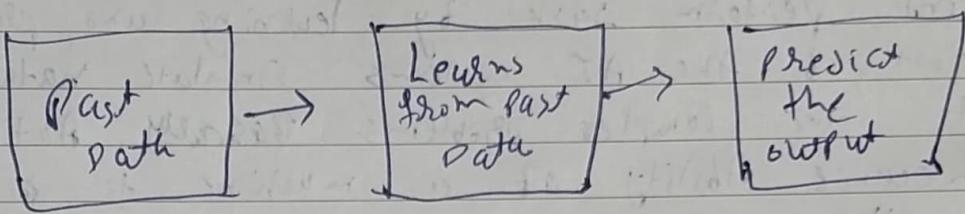
↳ Application of AI

- machine translate
- Self Driving Vehicles
- AI Robot
- Speech recognition applications like Siri or Google

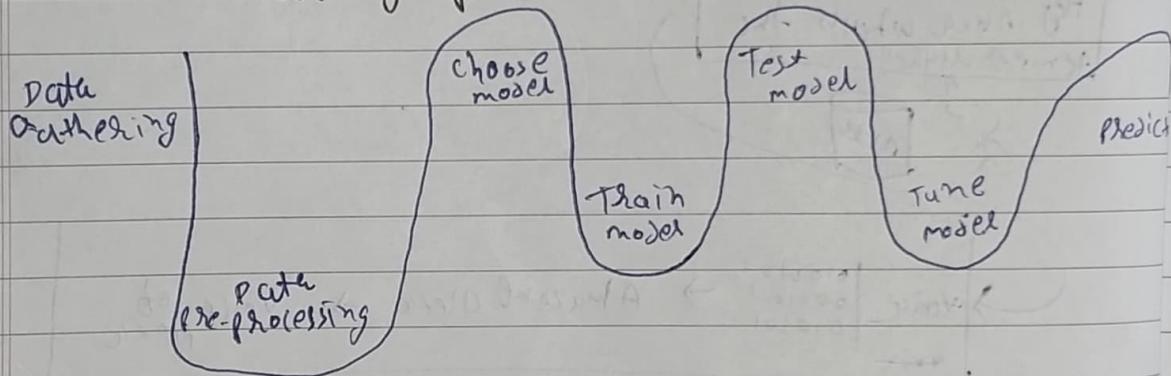
Machine Learning

Machine learning (ML) is a discipline of computer science that uses computer algorithms and analytics to build predictive models that can solve business problems.

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at task in T, as measured by P, improves with experience E.



Machine Learning processes



Application of ML

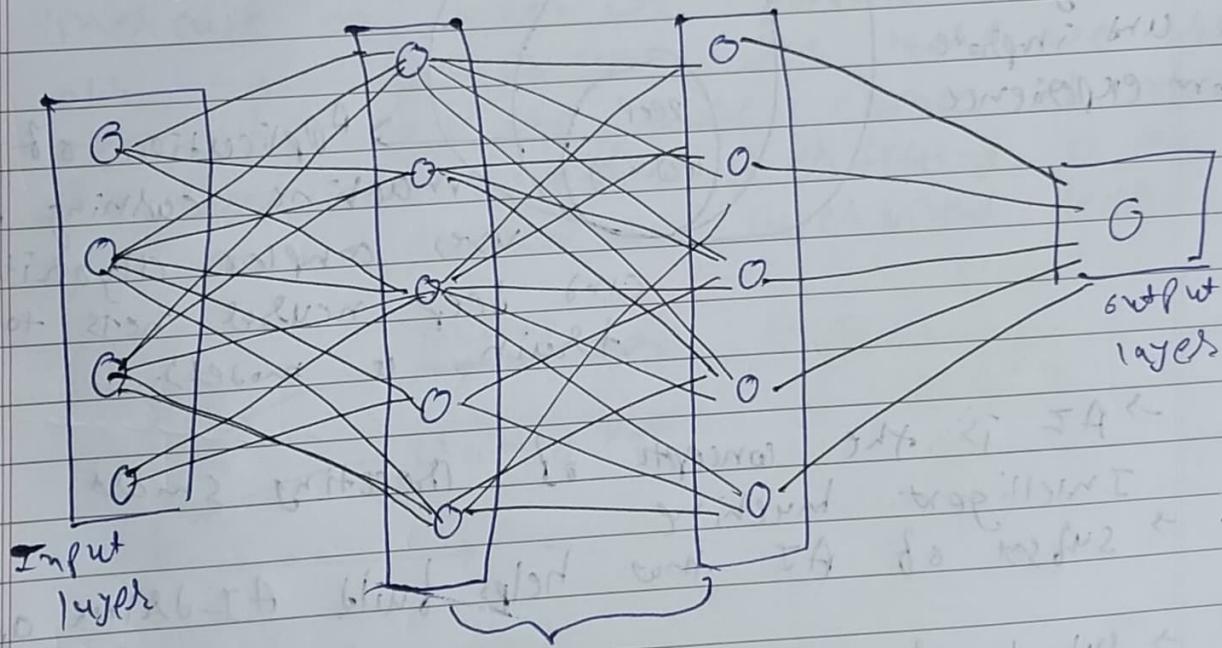
- Sales forecasting for different product
- Credit analysis in banking
- Product recommendation
- Stock price prediction

↳ Deep Learning

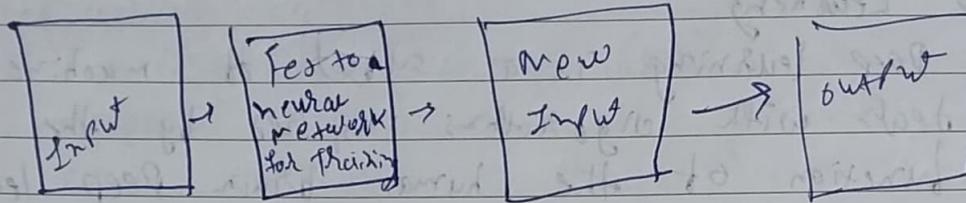
Deep learning is a subset of machine learning that deals with algorithms inspired by the structure and function of the human brain. Deep learning algorithms can work with an enormous amount of both structured and unstructured data. Deep learning's core concept lies in artificial neural networks, which enable machines to make decisions.

The major difference between deep learning vs machine learning is the way data is presented to the machine. ML algorithms usually require structures data, whereas deep learning network works on multiple layers of artificial neural networks.

↳ Simple neural network



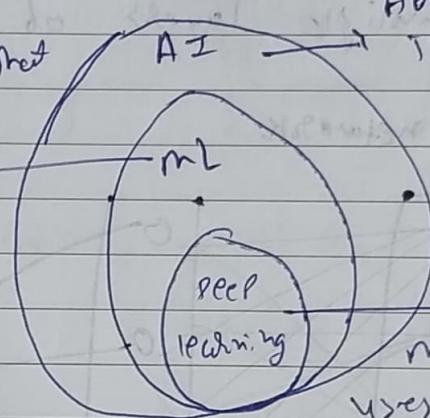
How It Works



↳ Application of Deep Learning

- cancer tumor detection
- captionbot for captioning an image
- music generating
- Image coloring
- object detection

Application of AI that
Allows a system to
automatically
learn and improve
from experience



Ability of a machine to
imitate intelligent human
behavior

→ Application of
machine learning that
uses complex algorithms
and deep neural nets to
train a model.

AI → AI is the concept of creating smart
intelligent machine

ML → subset of AI that helps build AI-driven applicati

Deep
learning → subset of ML that uses vast volumes of
data and complex algorithms to train a
model

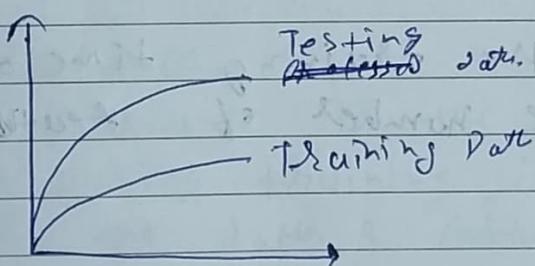
Q-2

Explain what is overfitting and underfitting.
 Explain the techniques to avoid overfitting and underfitting.

→ overfitting

overfitting occurs when our machine learning model tries to cover all data point or more than the required data points present in the given dataset. Because of this, the model starts picking noise and inaccurate values present in the dataset and all these factors reduce the efficiency and accuracy of the model. The overfitting model has low bias and high variance.

The chance of occurrence of overfitting increase as much we provide training to our model. It means the more we train our model, the more chance of causing the overfitted model. Overfitting is the most main problem in supervised learning.



Techniques to avoid

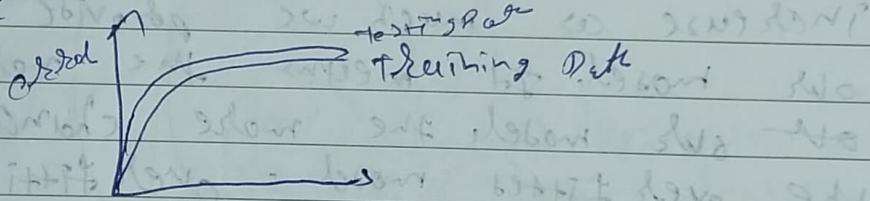
- Cross-validation
- Training with more Data
- Removing features
- Early stopping the training

Regulation
- Ensembling

Underfitting

- Underfitting occurs when our machine learning model is not able to capture the underlying trend of the data. To avoid the overfitting in the model, the fit or training data can be stopped at an early stage, due to which the model may not learn enough from the training data. As a result, it may fail to fit the best fit of the dominant trends in the data.

In the case of underfitting, the model is not able to learn enough from the training data, and hence it reduces the accuracy and produces unreliable predictions. Underfitting has high bias and low variance.



Technique to avoid

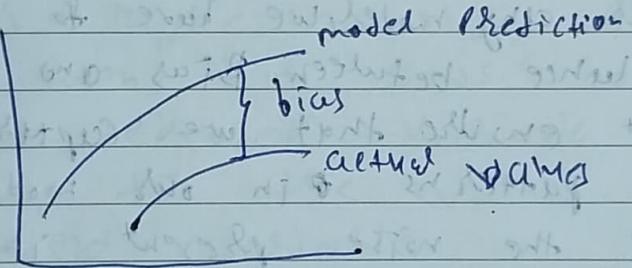
- Increasing the training time of model
- increasing the number of features.

Q.3 Explain the terms bias and variance with respect to machine learning. Discuss the bias-variance trade-off in machine learning.

↳ Bias

To make predictions, our model will analyze our data and find patterns in it. Using these patterns we can make generalizations about certain instances in our data. Our model after training learns these patterns and applies them to the test set to predict them.

Bias is the difference between our actual and predicted values. Bias is the simple assumptions that our model makes about our data to be predict new data.



↳ Variance

Variance is the very opposite of bias. During training, it allows our model to see the data a certain number of times to find patterns in it. If it does not work on the data for long enough, it will not find patterns and bias occurs. On the other hand, if our model is allowed to view the data too many times, it will learn very well for only the data. It will capture most patterns.

in the data, but it will also learn from the unnecessary data present or from the noise.

We can define variance as the model's sensitivity to fluctuations in the data. Our model learns from noise. This will cause our model to consider trivial features as important.

~~Bias~~ Variance

The variability of a model's prediction for a given data point which tells us the spread of our data is called the variance of the model.

~~Bias~~ Bias-variance tradeoff

For any model, we have to find the perfect balance between Bias and Variance. This just ensures that we capture the essential patterns in our model while ignoring the noise present in it. This is called Bias-variance Tradeoff. It helps optimize the error in our model & keeps it as low as possible.

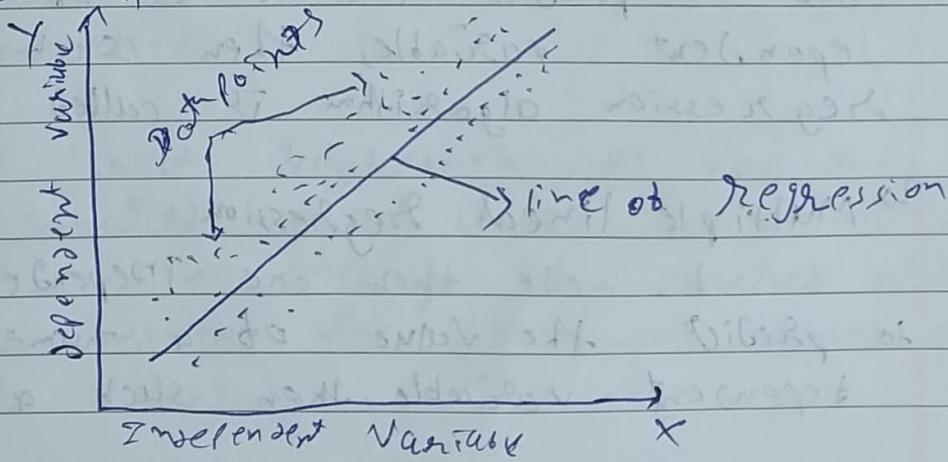
Ques.

Discuss in detail the application of linear regression for prediction

Linear regression is one of the most easiest and most popular machine learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/redd or numeric variable such as Sales, Salary, age, Product price, etc.

Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (x) variables, hence called as linear regression. Since linear regression shows the relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.

The linear regression model provides a slope straight line representing the best relationship between the variables & consider the



We can represent a linear regression as:

$$\underline{y = \beta_0 + \beta_1 x + \epsilon}$$

here,

y = Dependent Variable (Target Variable)

x = Independent Variable (Predictor Variable)

β_0 = Intercept of line

(Given as additional degree of freedom)

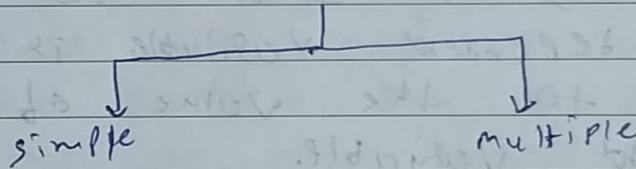
β_1 = Linear regression Coefficient

(Scale factor to each input value)

ϵ = random error

The value for x and y variable are training datasets for linear regression model representation.

Linear Regression



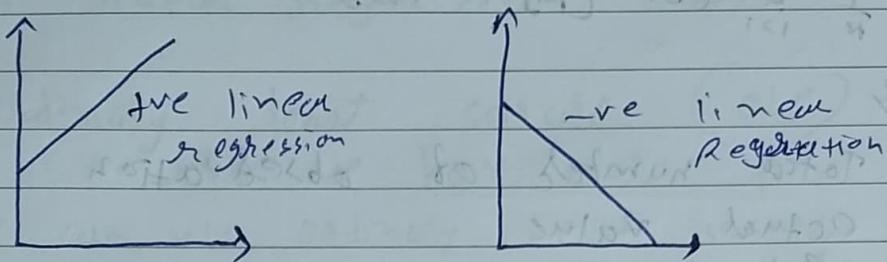
↳ Simple linear regression

If a single independent variable is used to predict the value of a numerical dependent variable, then such a linear regression algorithm is called simple.

↳ Multiple linear regression

If more than one independent variable is used to predict the value of a numerical dependent variable, then such a linear

Regression algorithm is called multiple linear regression



(cost function):

- The different values for weights or coefficients of line (P.w.) given the different line of regression, and the cost function is used to estimate the values of the coefficient for the best fit line.
- Cost function optimizes the regression coefficients or weights. It means measures how a linear regression model is performing.
- We can use the cost function to find the accuracy of the mapping function which maps the input variable to the output variable. This mapping function is also known as hypothesis function.

For linear regression, we use the mean squared error (MSE) cost function, which is the average of squared error occurred between the predicted values and actual values.

$f(x)$ can be written as

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - (\beta_1 x_i + \beta_0))^2$$

where,

N = Total number of observation

y_i = actual value

$(\beta_1 x_i + \beta_0)$ = predictive value.

Residuals:-

The distance between the actual value and predicted values is called residual. If the observed points are far from the regression line, then the residual will be high, and so the cost function will be high. If the scattered point are close to the line, then the residual will be small and hence the cost function will be small.

a.5 Briefly describe any two techniques used to compute the parameters of the linear regression model

→ (i) Ordinary least square (OLS)

when we have more than one input we can use ordinary least squares to estimate the values of the coefficients. The OLS procedure seeks to minimize the sum of the square error residuals. This means that given a regression line through the data we calculate the distance from each data point to the regression line, square it, and sum all of the square errors (SSE) together. This is the quality that OLS seeks to minimize.

$$SSE = \frac{\sum_{i=1}^n d_i^2}{n}$$

This approach treats the data as a matrix and uses linear algebra operations to estimate the optimal values for the coefficients. It means that all of the data must be available and you must have enough memory to fit the data and perform matrix operations. It is unusual to implement the OLS procedure yourself unless as an exercise in linear algebra. It is more likely that you will call a procedure in a linear algebra

library. This procedure is very fast to calculate.

$$Y = \beta_0 + \beta_1 X$$

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\beta_0 = \bar{Y} - \beta_1 \bar{x}$$

where,

x_i = independent variables

\bar{x} = average of independent variable

Y = dependent variables

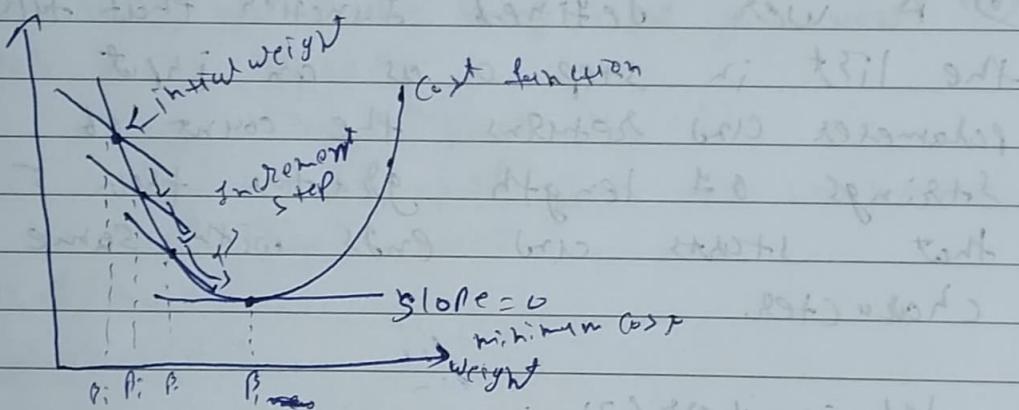
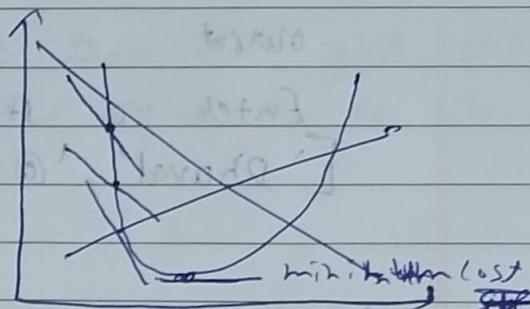
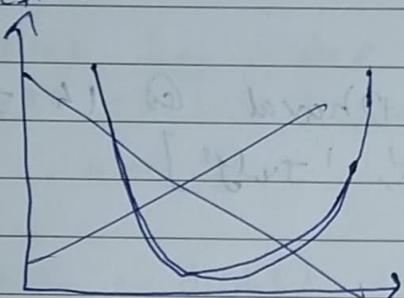
\bar{Y} = average of dependent variable

(ii) Gradient Descent

Gradient Descent is known as one of the most commonly used optimization algorithm to train machine learning models by means of minimizing errors between actual and expected results. Further, gradient descent is also used to train neural network.

In mathematical terminology, Optimization Algorithm refers to task of minimizing/maximizing an objective function $f(x)$ parameterized by x . Similarly, in ML, optimization is the task of minimizing the cost function parameterized by the

model's parameters. The main objective of gradient descent is to minimize the convex function using iteration of parameter updates once these machine learning models are optimized, these models can be used as your sub tools for AI and various computer science applications. It helps in finding the local minimum of a function.



$$J(\beta_0, \beta_1) = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

$$\hat{y} = \beta_0 + \beta_1 x_1 + g$$

$$\beta_i = \beta_i - \eta \frac{\partial J(\beta_0, \beta_1)}{\partial \beta_i}$$

Q.6

write a python program that does has the following:

- a) takes string inputs from the user and store in a list

```
S+r = list("Enter a string: ").split(" ")
print(S+r)
```

Output

```
Enter a string: Dharav@14 July
['Dharav', '@', '14', 'July']
```

- b) a user defined function that takes the list in step a as an input parameter and returns the count of strings of length greater than 5 that starts and ends with same character

→ def inputstr():

```
list = list(input("Enter your string: ").split(" "))
print(list)
```

count = 0

for x in list:

if len(x) > 5:

if x[0] == x[len(x)-1]:

count = count + 1

list.clear()

return count

input str()

Output:

Enter your string: h a l b d j h a d u g j s d h k a l e k n o i a u k
 (or any other string)

or 4

Q.7 write a Python program that does / has the following

a.) Takes key value inputs from the user and store in a dictionary - The value is either of numbers.

```
mydisk = dict()
n = int(input("How many input you want:"))
for x in range(n):
    key = input("Enter Key:")
    value = int(input("Enter value in numeric form:"))
    mydisk[key] = value
print(mydisk)
mydisk.clear()
```

Output:

How many input you want: 3

Enter Key : K1

Enter Value in numeric form: 1

Enter Key : K2

Enter Value in numeric form: 2

Enter Key : K3

Enter Value in numeric form: 3

{'K1': 1, 'K2': 2, 'K3': 3}

b.) A user defined function that takes the dictionary in step a) as input parameter and return the key that has minimum number of unique elements in its list

Q.8 A box contain 3 blue marble, 4 red, 6 green marble and 2 yellow marbles. If three marbles are picked at random, what is the probability that they are all blue?

Given that there are

3 blue marble

4 red marble

6 green marble

2 yellow marble

$$\text{Total marble} = 15$$

Probability that they all be blue is

$$P(\text{all are blue}) = \frac{3C_3}{15C_3}$$

$$= \frac{1}{45}$$

Q.9 A bag contains 6 blue balls and 4 red balls. 3 balls are picked at random. What is the probability that none of them is Red?

Given that there are

6 Blue Balls } total balls = 10
4 Red Balls

Probability of getting Red 1/2 times

$$P(\text{Getting Red}) = \frac{6C_3}{10C_3}$$

$$\frac{6 \times 5 \times 4}{10 \times 9 \times 8} \times 2$$

$$= \frac{1}{30}$$

Here for we can say none of them

is Red is equal all are blue

$\therefore P(\text{Getting none of them is Red})$

$$= P(\text{Getting blue}) = \frac{6C_3}{10C_3}$$

3

$$= \frac{6 \times 5 \times 4}{10 \times 9 \times 8} \times 2$$

$$= \frac{1}{8}$$

Q.10 The probability that A speaks truth is $\frac{3}{5}$ and that of B speaking truth is $\frac{4}{7}$. What is probability that they agree in stating the same fact?

$$\rightarrow \text{Given } P(A) = \frac{3}{5} \text{ speaking Truth}$$

$$P(B) = \frac{4}{7} \text{ speaking Truth}$$

Find probability that A and B agree to the same fact.

\rightarrow Since A and B are independent,

$$\therefore P(A \cap B)$$

$$= P(A) + P(B) \quad [\text{speaking of both Truth}]$$

$$= \frac{3}{5} \times \frac{4}{7}$$

$$= \frac{12}{35}$$

\rightarrow Probability both speaking false

$$= [1 - P(A)] + [1 - P(B)]$$

$$= [1 - \frac{3}{5}] \times [1 - \frac{4}{7}]$$

$$= \frac{2}{5} + \frac{3}{7}$$

$$= \frac{6}{35}$$

↪ Probability, both agree the same

$$= P(\text{Speaking Truth}) + P(\text{Speaking False})$$

$$= \frac{12}{35} + \frac{6}{35}$$

$$= \frac{12+6}{35}$$

$$\boxed{= \frac{18}{35}}$$

Q.12 Explain in detail any five applications of machine learning.

1) Virtual personal Assistant:

We have various virtual assistants such as Google Assistant, Alexa, Cortana, Siri, As the name suggests, they help us in finding the information using our voice instruction. These assistants can help us in various ways just by our voice instructions such as play music, call someone, open an email, scheduling an appointment etc.

These assistants record our voice instruction, send it over the server on a cloud, and decode it using ML algorithms and act accordingly.

2) Transportation and commuting (cabs) (Uber)

Online cab booking using ML to an extent. It provides a personalized application which is unique to you. Automatically detects your location and provides options to either go home or office or any other frequent places based on your history and patterns.

It uses ML algorithm layered on top of historical trip data to make a more accurate ETA prediction. With the implementation of ML, they saw a 26% accuracy in delivery and pickup.

3) Self Driving Cars (TESLA)

ML plays a very important role in self driving cars and we all heard about TESLA, the leader in this business and their current AI is driven by hardware manufacturer NVIDIA, which is based on unsupervised learning algorithm.

4) Traffic Alerts (maps)

Google maps use whenever we go to some new place and it shows us which route may have traffic which is more faster. But how does it know that?

It's a combination of people currently using the service, historic data of that route collected over time and few tricks acquired from other companies. Everyone using maps is providing their location, average speed, the route in which they are travelling which in turn help Google collect massive data about traffic, which makes them predict the upcoming traffic and adjust your route according to it.

5) Email spam and malware filtering

Whenever we receive a new mail, it is filtered automatically as important, normal, and spam. We always receive an important mail in our inbox with

the important symbol and spam emails in our spam box and the technology behind this is machine learning. Below are some spam filters used by Gmail:

- content filter
- header filter
- general blacklist filter
- rules-based filters
- permission filter

Some ML algorithms such as multi-layered perceptron, decision tree and naive Bayes classifiers are used for email spam filtering and malware detection.

Q.17 Differentiate between Pearson's correlation and Spearman's correlation. Give examples illustrating the use of both the measures.

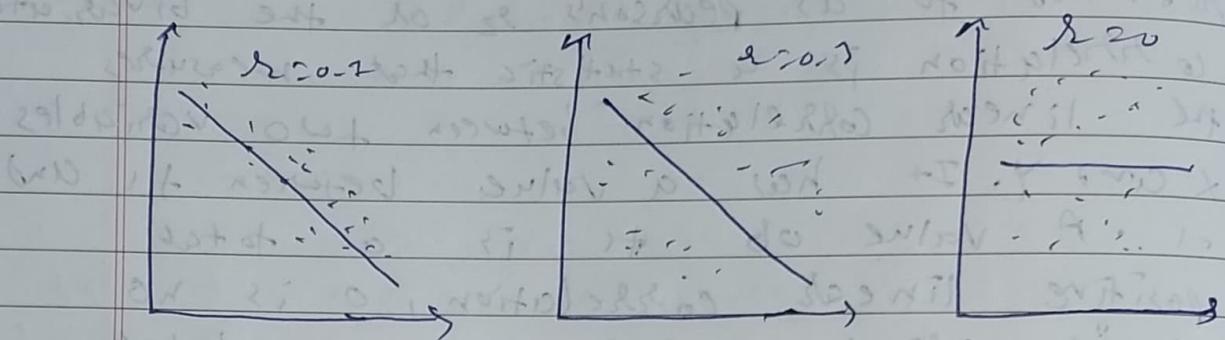
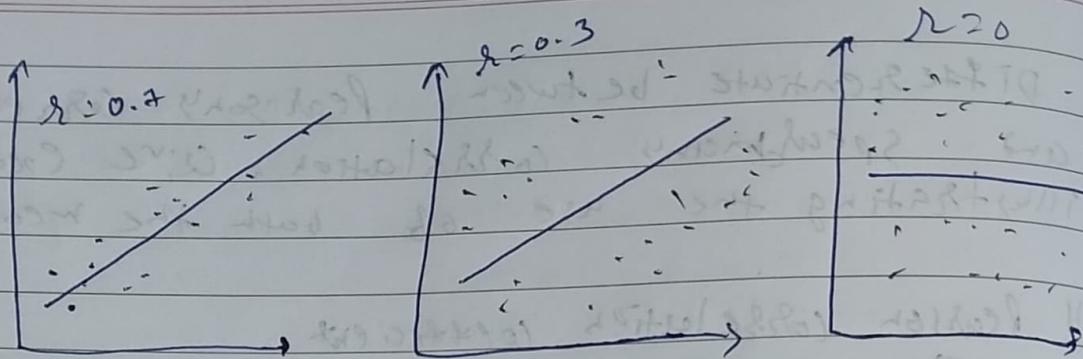
→ 1) Pearson correlation coefficient:

The Pearson correlation coefficient also referred to as Pearson's r or the bivariate correlation is a statistic that measures the linear correlation between two variables X and Y . It has a value between $+1$ and -1 . A value of $+1$ is a total positive linear correlation, 0 is no linear correlation, and -1 is a total negative linear correlation.

The Pearson correlation can evaluate only a linear relationship between two continuous variable (A relationship is linear only when a change in one variable is associate with a proportional change in the other variable)

Example: we can use the Pearson correlation to evaluate whether an increase in age leads to an increase in blood pressure.

→ Below is an example of how the Pearson correlation coefficient (r) varies with the strength and the direction of the relationship between the two variables. Note that when no linear relationship could be established, the Pearson coefficient yields a value of zero.



\hookrightarrow Spearman's Correlation Coefficient

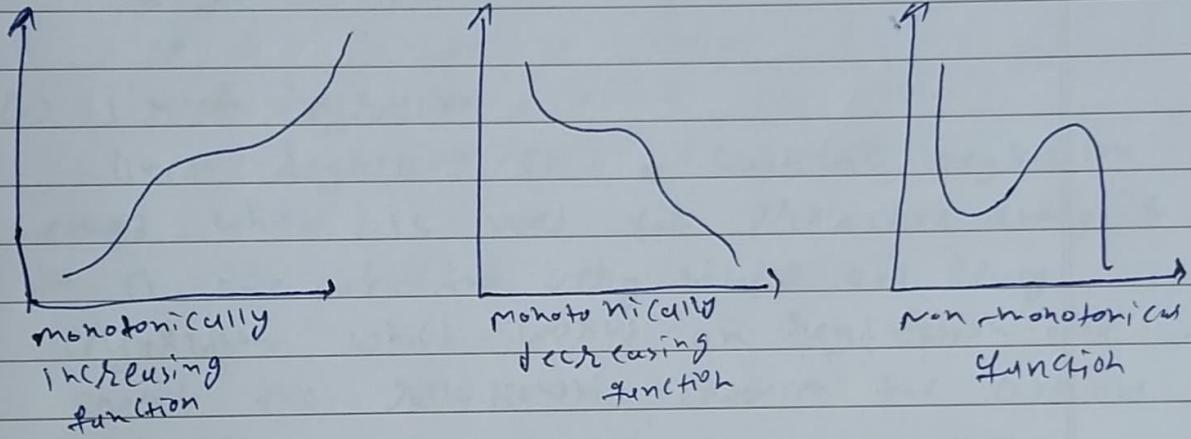
Spearman's rank correlation coefficient or Spearman's ρ (named after Charles Spearman) is a nonparametric measure of rank correlation (statistical dependence between the ranking of two variables.) It assesses how well the relationship between two variables can be described using a monotonic function.

After word The Spearman correlation can evaluate a monotonic relationship between two variables - continuous or ordinal and it is based on the ranked values for each variable rather than the raw data.

What is monotonic relationship?

- A monotonic relationship is a relationship that does one of the following:
 - as the value of one variable increases, so does the value of the other variable, or,
 - as the value of one variable increases, the other variable value decreases.

BUT, not exactly at a constant rate whereas in a linear relationship the rate of increase/decrease is constant



Example:- whether the order in which employees complete a test excel size is related to the number of month they have been employed or correlation between the IQ of a person with the number of hours spent in front of TV per week.

Q.1 Explain in brief any five supervised learning models.

Supervised learning can be divided into two types of problems:

Regression

Classification

- Linear regression → Random forest
- Non-linear regression → Decision trees
- Bayesian linear regression → Logistic regression
- Polynomial Regression → Support Vector machines
- Regression Trees

(i) Linear regression

- Linear regression is a statistical regression method which is used for predictive analysis.

- It is one of the very simple and easy algorithms which works on regression and shows the relationship between the outcome variable.

- It is used for solving the regression problems in machine learning.

- Linear regression shows the linear relationship between the independent variable (x) and dependent variable (y), hence called linear regression.

- If there is only one input variable given such linear regression called simple linear regression and if there

If we have more than one input variable, then this is called multiple linear regression.

→ Equation of linear regression

$$Y = mX + C$$

Here,

Y = dependent variable (target variable)

X = independent variable (predicted variable)

m & C = linear coefficient

(ii) Non-linear regression

- non-linear regression is a form of regression analysis in which observational data are modeled by a function which is a non-linear combination of the model parameters and depends on one or more independent variables. The data are fitted by a method of successive approximations.

A simple non-linear regression model is expressed as follows:

$$Y = f(X, \beta) + \epsilon$$

where

X = vector of predictor

β = vector of k parameters

$f(\cdot)$ = known regression function

ϵ = error term

Alternatively, the model can also

be written as

$$Y_i = h[X_i^{(1)}, X_i^{(2)}, \dots, X_i^{(n)}; \theta_1, \theta_2, \dots, \theta_p] + \epsilon_i$$

Where,

y = responsive variable

h_2 = function

x = input data matrix (n × p)

θ = parameters to be estimated

Example

k-nearest Neighbours

kernel SVM

Naive Bayes

Decision Tree

(iii) Naive Bayes algorithm

Naive Bayes algorithm is a type of supervised

learning algorithm, which is

based on Bayes theorem and used for solving classification problems.

This is mainly used in text classification that includes a high-dimensional training set.

Bayes theorem

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

(iv) Logistic Regression

- logistic Regression is another supervised learning algorithm which is used to solve the classification problem. In classification problem, we have dependent variable in a binary or discrete format such as 0 or 1.

- logistic regression algorithm which work on the concept of probability.

- logistic regression is a type of regression, but it is different from the linear regression algorithm in the term how they are used.

- logistic regression uses Sigmoid function or logistic function which is a complete cost function. This sigmoid function is used to model the data in logistic regression. The function can be represented as

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

where

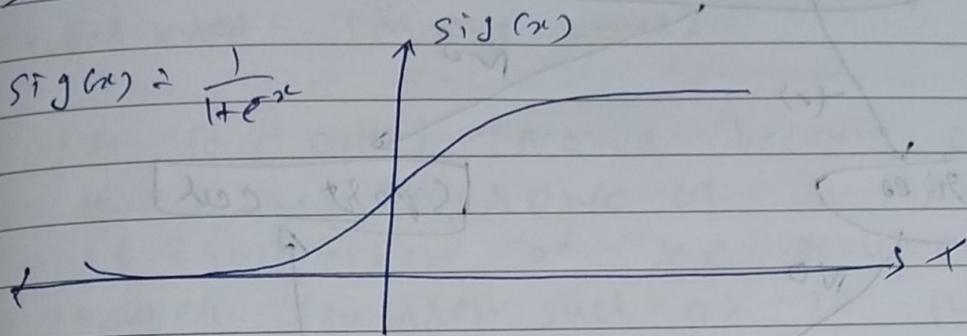
$f(x)$ = output between the 0 and 1.

x = input to the function

e = base of natural logarithm

when we provide the input values to the function, it gives the

S-curve as follows.



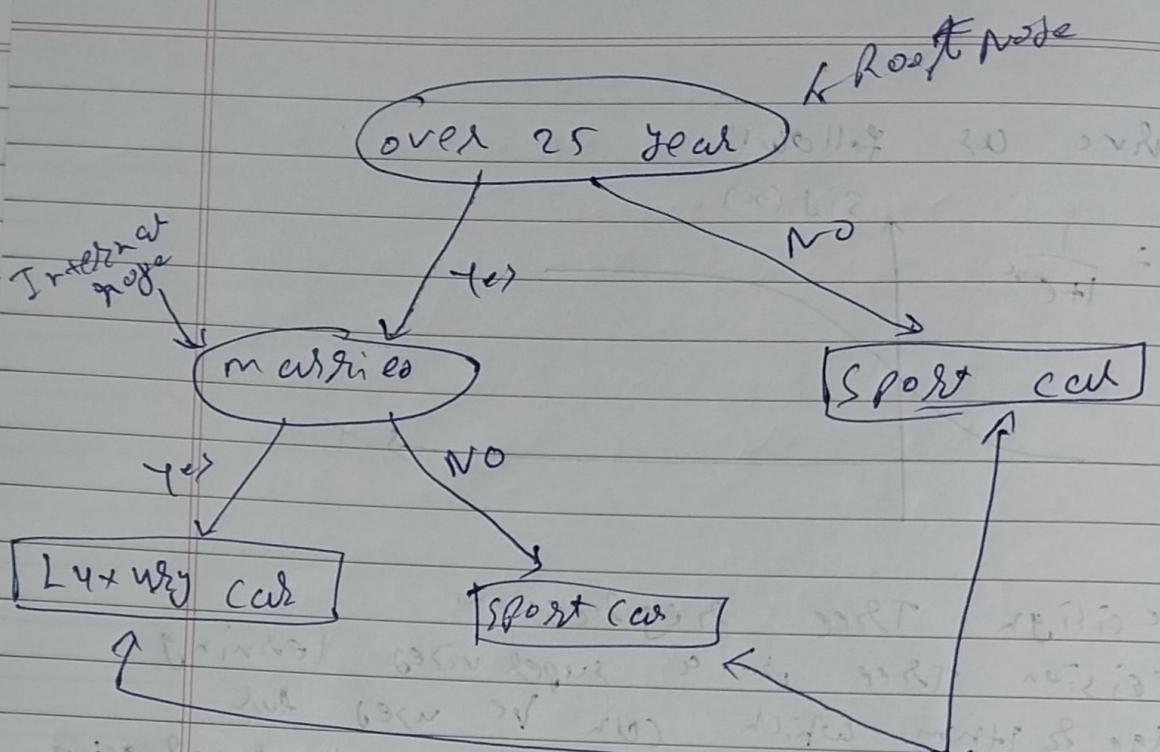
(v) Decision Tree regression

- Decision tree is a supervised learning algorithm which can be used to solve both classification and regression problem

- It can solve problems for both categorical and numerical data

- Decision tree regression builds a tree-like structure in which each internal node represents the "test" for an attribute such branch represent the result of the test, and each leaf node represents the final decision or result

- A decision tree is constructed starting from the root node/purest node (dataset), which splits into left and right child nodes (subset of dataset). These child nodes are further divided into their child nodes, and themselves become the parent node of those nodes. Consider the following



a.15 what is Naive in Naive Bayes supervised learning model?

→ It is called Naive because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. such as if the fruits is identified on the base of color, shape, and taste then red, spherical and sweet fruit is recognized as an apple. hence each feature individually contribute to identify that it is an apple without depending on each other.

a.16 Explain the terms loss function and cost function. Differentiate between both.

The loss function is a method of evaluating how well your machine learning algorithm models your features dataset. In other words loss function are a measurement of how good your model is in term of predicting the expected output.

Cost Function

Cost function measures the performance of a machine learning model for a data set. Cost function quantifies the error between predicted and expected values and presents that error in the form of a single real number.

Depending on the problem, cost function can be formulated in many different ways.

Q-17) Differentiate between parametric model and non-parametric model

Parametric Model

In case of parametric models, the assumption related to the function form is more linear model is considered. In case of non-parametric models, the assumption about the functional form is not made.

Parametric model is much easier to fit than non-parametric models because parametric model machine learning models only require the estimation of a set of parameters as the model is identified prior as linear model. In case of non-parametric models, one needs to

estimate some arbitrary function which is a much difficult task.

- parametric models often do not match the unknown function we are trying to estimate. The model performance is comparatively lower than non-parametric model.

The estimates done by the parametric models will be further from being true.

- Parametric models are interpretable unlike the non-parametric model. This essentially means that one can go for parametric models when the goal is to find inference instead, one can go for non-parametric model when the goal is to make prediction with higher accuracy and interpretability of inference is not the key ask.

Q18 Define the cost function for logistic regression model

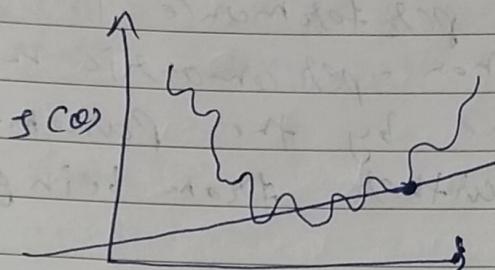
→ In the case of linear Regression, the cost function is -

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} [h_\theta(x^{(i)}) - y^{(i)}]^2$$

But for logistic Regression,

$$h_\theta(x) = g(\theta^T x)$$

It will result in a non-convex cost function. But this result in cost function with local optima's which is a very big problem for gradient descent to compute the global optima.

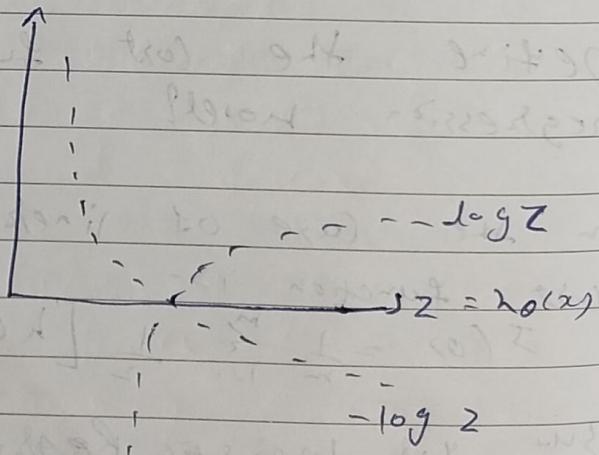
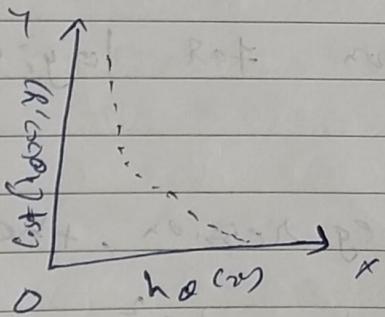


Non convex Cost Function

So, for logistic regression the cost function is

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)), & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

if $y=1$



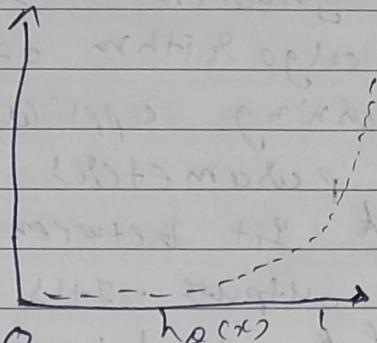
$$(o_0 + \theta_1 x) = 0 \quad \text{if } y = 1, h_{\theta}(x) = 1$$

But as

$$h_{\theta}(x) \rightarrow 0$$

$(o_0 + \theta_1 x) \rightarrow \text{Infinity}$

If $y = 0$



so,

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} 0, & \text{if } h_{\theta}(x) = y \\ \infty, & \text{if } y = 0 \text{ and } h_{\theta}(x) \rightarrow \infty \\ \infty, & \text{if } y = 1 \text{ and } h_{\theta}(x) \rightarrow 0 \end{cases}$$

$$(o_0 + \theta_1 x, y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

To fit the parameter θ , $J(\theta)$ has to be minimized and for the gradient descent it is required.

Gradient Descent - looks similar to that of linear regression but the difference lies in the hypothesis $h_{\theta}(x)$.

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) z_j^{(i)}$$

Q.19

with respect to gradient descent
define the following terms

→ a) Stochastic gradient descent:

Stochastic gradient descent is an optimization algorithm often used in machine learning applications to find the model parameters that correspond to the best fit between predicted and actual outputs. It's an inexact but powerful technique. Stochastic gradient descent is widely used in machine learning applications.

b) Batch gradient descent:

In batch gradient descent, all the training data is taken into consideration to take a single step. We take the average of the gradients of all the training examples and then use that mean gradient to update our parameters. So that's just one step of gradient descent in one epoch.

c) mini-batch gradient descent →

We use a batch of a fixed number of training examples which is less than the actual dataset and called it a mini-batch. Doing this helps us archive the advantages of both the former variations we

So, after creating the mini-batch of fixed size, we do the following steps in one epoch:

1) Pick a mini batch

2) Feed it to Neural Network

3) calculate the mean gradient of the mini-batch

4) use the mean gradient we calculated in step 3 to update the weights

5) Repeat step 1 to 4 for the mini-batches we created

Q.26 What is Regularization? Explain any two regularization techniques.

→ Regularization :-

Regularization refers to a set of different techniques that allow lower the complexity of a neural network model during training, thus prevent the overfitting

L_2 & L_1 regularization

L_1 and L_2 are the most common types of regularization. These update the general cost function by adding another term known as the regularization term

Cost function = Loss (say, binary cross entropy) + Regularization term

Due to addition of this regularization term the value of weight matrices decrease because it assumes that a learned network with smaller weight matrices leads to simpler models.

Therefore, it will also reduce overfitting to quite an extent.

However, this regularization term differs in L1 and L2. In L2, we have

$$\text{Cost function} = \text{Loss} + \frac{\lambda}{2m} \|w\|^2$$

Here, λ is the regularization parameter. It is the hyperparameter whose value is optimized for better result. L2 regularization is also known as weight decay as it forces the weight to decay towards zero (but not exactly zero).

In L1, we have

$$\text{Cost function} = \text{Loss} + \frac{\lambda}{2m} + \frac{1}{2m} \|w\|_1$$

In this, we penalize the absolute value of the weights. Unlike L2, the weights may be reduced to zero.

Here, hence if it is very useful when we are trying to compress our model. otherwise, we