

## Experiment No. 1

Aim: Select an appropriate dataset and perform the following preprocessing steps

i) Imputation : it is the process of replacing missing values in the dataset there are several techniques to handle missing values

- a) Mean/Median/Mode : Imputation replace missing values with the mean, median, mode of the respective column.
- b) Forward and Backward fill : Fill missing values using the previous or next observation.
- c) Interpolation use linear or polynomial interpolation to estimate missing values.
- d) K-nearest Neighbours use the CNN algorithm to impute missing values based on the values of their nearest neighbour

2) Anomaly Detection: It involves identifying or abnormal data points in the datasets.

- a) Statistical Methods
- b) Machine learning models
- c) time series methods

3) Standardization :- Transforms data to have a mean of standard deviation of one useful for algorithms that assume normally distributed data

4) Normalization:- Transform data to have a mean of zero & standard deviations of one useful for algorithms that assume normally, which is useful for algorithms that rely on distance measurement.

5) Encoding : used to convert categorical data into numerical form

a) Label encoding

b) One-hot encoding

c) Ordinal encoding

d) Frequency encoding

(A) ~~Label~~  
Frequency

## Experiment No. 2

Aim: Implement Gradient Descent algorithm to minimize a given function.

Theory: Gradient Descent is an optimization algorithm used to minimize a given objective function particularly in machine learning & optimization problems. It is an iterative approach that seeks to find the minimum of a function by updating parameters in the direction

1) Objective function:-  
the objective function, also known as the cost function or loss function measures how good or bad a set of parenthesis is. The goal of descent is to find the set of parameters that minimize this function.

2) Initialization:- You start by initializing the parameters of the model with some initial values.

3) Gradient Calculation:- Calculate the gradient of the objective function with some initial values.  
The gradient is a vector that points in a direction of the increases of the function.  
In the context of optimization you need to move in the opposite direction to reach the minimum.

4) Update parameter : Update the parameters by moving them a small step in the direction of the negative gradient this step in the direction of the negative gradient this step size is determined by a parameter called the learning rate.

5) Iteration : Repeat step 3 & 4 until a stopping criterion is met. The criterion would be maximum number of iterations, the change in the objective function very small or some other measure indicating convergence.

(A) Stochastic

## Experiment No. 3

Aim:- Implement Simple Linear Regression

Theory :-

Linear Regression model studies the relationship between a single dependent variable ' $y$ ' & one or more independent variables ' $x$ ' using a best fit straight line function. This line can be characterized by the slope intercept with one of the axis. Hence given data points goal is to find a line which indicates finding slope & intercept.

$$y = \beta_0 + \beta_1 x + \epsilon$$

Where  $\epsilon$  - random error associated with data point so we assume underlying function from which the data is generated  $y = \beta_0 + \beta_1 x + \epsilon$ . the given data points, goal is to find out  $\beta_0$  &  $\beta_1$ .

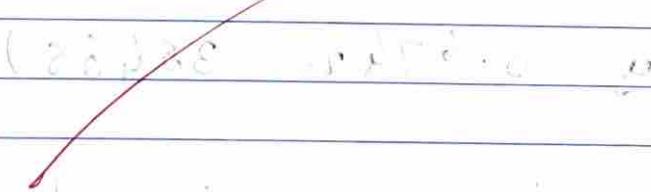
$$y = \beta_0 + \beta_1 x + \epsilon$$

where  $\beta_0$  = population of  $y$  - intercept

$\beta_1$  =  $\beta_0$  population slope

$\epsilon$  = random error

e.g.



weight	weight	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$C_{xi,yi}$
151	6.3	2.8	-2.3	6.4	7.84
144	81	20.2	20.2	317.14	607.04
138	56	-15.8	-15.8	146.48	249.67
186	91	32.2	32.2	827	1036.8
128	47	25.8	-25.8	-672	656.6
136	57	-17.8	-17.8	147	116
149	76	25.2	25.2	616	84.64
163	72	9.2	9.2	84	84.64
152	62	-1.8	-1.8	44	3.24
131	48	22.8	-22.8	264	3127.6
153	65	8	8	64	

for m<sub>0</sub> :  $b_1 = 2649.6$

$= 3927.6$

Value of  $b_0 = 0.67436$

③  $b_0 = \frac{1}{n} (\sum y_i - b_1 \sum x_i) = \bar{y} - b_1 \bar{x}$

$$\begin{aligned}
&= 0.65 - 3(0.67436 \times 183.8) \\
&= -38.4551 \\
&= 0.67436 - 38.4551
\end{aligned}$$

$y = 0.67436 - 38.4551$

A<sup>x</sup> C<sup>y</sup>  
S<sup>z</sup>

## Experiment No. 4

Aim :- Implement Logistic Regression using machine learning method without using sklearn

Theory :- Logistic Regression estimates the probability of an event occurring, such as 'voted or didn't vote', based on a given dataset of independent variables since the outcome is a prob the dependent variable is bounded between 0 & 1

1] Binomial : - In this there can be 2 or more possible unordered types of the dependent variables such as Yes, No.

2] Multinomial : In this there can be 3 or more possible ordered type of the dependent variables such as cat, dogs, sheep, etc.

3] Ordinal: In ordinal logistic regression there can be 3 or more possible ordered types of dependent variables such as low, medium or higher assumption.

4] Independent Observations : Each observations is depend of the others meaning there is no correlation b/w input variables.

5] Binary dependent variable It takes the assumption that the dependent variable must be binary or discrete

3) Large sampling :- the sample size is large.

4) Linearity relationship b/w dependent variables.

Let the independent input be

$$\alpha = \begin{bmatrix} x_{11} & \dots & x_{1m} \\ x_{21} & \dots & x_{2m} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nm} \end{bmatrix}$$

& the dependent variable is having only binary values.

$$= 1 \text{ if class } 1$$

$$\text{Probability } P(y=1) = \sigma(z)$$

Final logistic regression

$$P(w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n)$$

(A) Sigmoid function

## Experiment No. 5

**Aim :-** Implement Decision Tree classification Algo\*  
**Theory :-** A decision tree is a tree-like structure that represents a set of decisions and their possible consequences. Each node in the tree represents a decision. The leaves of the tree represent the final decisions. The leaves of the tree represent the final decisions.

Consider the dataset

age	competition	Types	Profit
Old	yes	Software	down
Old	no	software	down
Old	No	Hardware	down
Mid	yes	Software	down
Mid	yes	Hardware	down
Mid	No	Hardware	up
Mid	No	Software	up
New	Yes	Software	up
New	No	Hardware	up
New	No	Software	up

$$\text{Gini}(7) = \left[ 1 - \left[ \frac{5}{10} + \left( \frac{5}{10} \right)^2 \right] \right] = 0.5$$

gini impurity

Age: old node, mid node, new node

$$\text{split t} = \frac{3}{10} \text{ gini old} + \frac{4}{10} \text{ gini mid} + \frac{3}{10} \text{ gini new}$$

$$= \frac{3}{10} \left[ 1 - \left( \frac{0}{3} \right)^2 + \left( \frac{2}{3} \right)^2 \right] + \frac{4}{10} \left[ 1 - \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right] + \frac{3}{10} \left[ 1 - \left( \frac{1}{3} \right)^2 \right]$$

$$= 0.2$$

### \* Competition

$$\text{split t} = \frac{4}{10} \text{ gini yes} + \frac{6}{10} \text{ gini no}$$

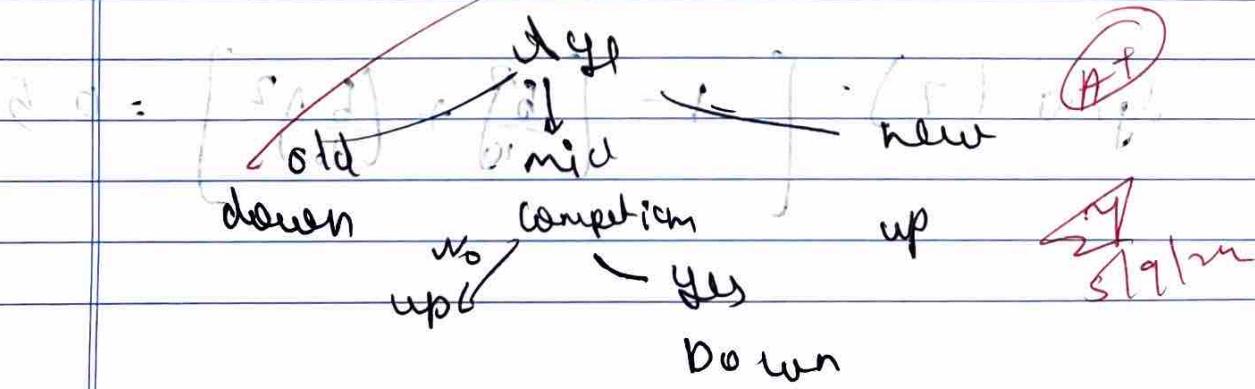
$$= 0.62$$

### \* type

$$\text{split t} = \frac{6}{10} \text{ gini (software)} + \frac{4}{10} \text{ gini (hardware)}$$

$$= 0.5$$

- Split values of age is the smallest age at root node



## Experiment 6

Aim: Implement SVM for Classification

Theory: SVM is a powerful Machine learning algorithm used for linear or non-linear classification, regression and even outlier detection tasks. SVM can be used for a variety of tasks such as text, such classification, image classification, spam detection, handwriting identification, gene expression analysis, face detection tasks. SVMs are adaptive & efficient in a variety of applications because they can manage high-dimensional and difficult non-linear relationships. SVM algorithms are very effective as we try to find the maximum separating hyperplane between the different classes accessible in the target feature space.

SVM terminology :-

Hyperplane: Hyperplane is the decision boundary that is used to separate the data points of different classes in a feature space.

Support Vectors: Support vectors are the closest data points to the hyperplane, which makes a critical job in deciding the hyperplane & margin.

Margin: Margin is the distance b/w the support vector & hyperplane. Hard margin: the hard margin hyperplane is a hyperplane that properly separates the data points of

of different categories without any misclassification. Soft margin: When the data is not perfectly separable or contains outliers, this SVM permits a soft margin. Each data point has a slack variable introduced by the soft margin SVM formulation, which softens the strict margin boundary & permits certain misclassification.

Conclusion:

We successfully implemented SVM for classification.

(X)

A  
Info/m

## Experiment 7

Aim :- Implement ensemble methods to combine different models.

Theory :

Ensemble methods are powerful techniques in ML that combine the predictions of multiple individual models to improve overall predictive performance.

Advantages:

- (1) Reduces Overfitting
- (2) Increases Model Robustness
- (3) Enhances prediction accuracy

(1) Voting ensemble:

Hard voting :- Multiple base models make predictions, final decision is the one with majority vote.

Soft voting : Similar to hard voting, but instead of selecting the class with most votes, it calculates weighted avg of class probabilities predicted by the base models

(2) Bagging :- Bagging aims to reduce the variance of an individual base model by training multiple base models on random subsets of the training data with replacement. The final prediction is often made by averaging or majority voting.

### (9) Boosting :-

They focus on improving the accuracy of weak base models sequentially. Each base model is trained to correct the mistakes made by the previous model.

### (10) Stacking :-

It combines multiple base models by training a meta-model on their predictions. The base models output serve as features for the meta model. It is known for its high predictive performance.

### (11) Ensemble Selection :-

These techniques aim to reduce the complexity of an ensemble by selecting a subset of base models that contribute the most to the ensemble's performance. Reduces computational overhead.

Other methods: Stacking with cross-validation, weighted, or voting ensembles, adaptive ensemble.

Conclusion :

Ensemble methods implemented

A. S.  
14/10/2021

## Experiment 8

Ques :- Implement Principal Component Analysis technique

Theory :-

PCA is a dimensionality reduction technique commonly used to transform high-dimensional data into a lower dimensional representation while preserving the most important information.

PCA : This by finding orthogonal axes along which the data exhibits the maximum variance.

These principal components are linear combination of the original features, allowing you to reduce the dimensionality of the data.

The steps & concepts of PCA are :-

- ① Centering the data : By subtracting the means of each feature from the corresponding data point
- ② Covariance matrix :
- ③ Eigen value decomposition
- ④ Eigen values & Eigen vectors

## (8) Selecting principal components

### (a) Projection:

Conclusion: Principal component analysis

is implemented

(A)

(B)

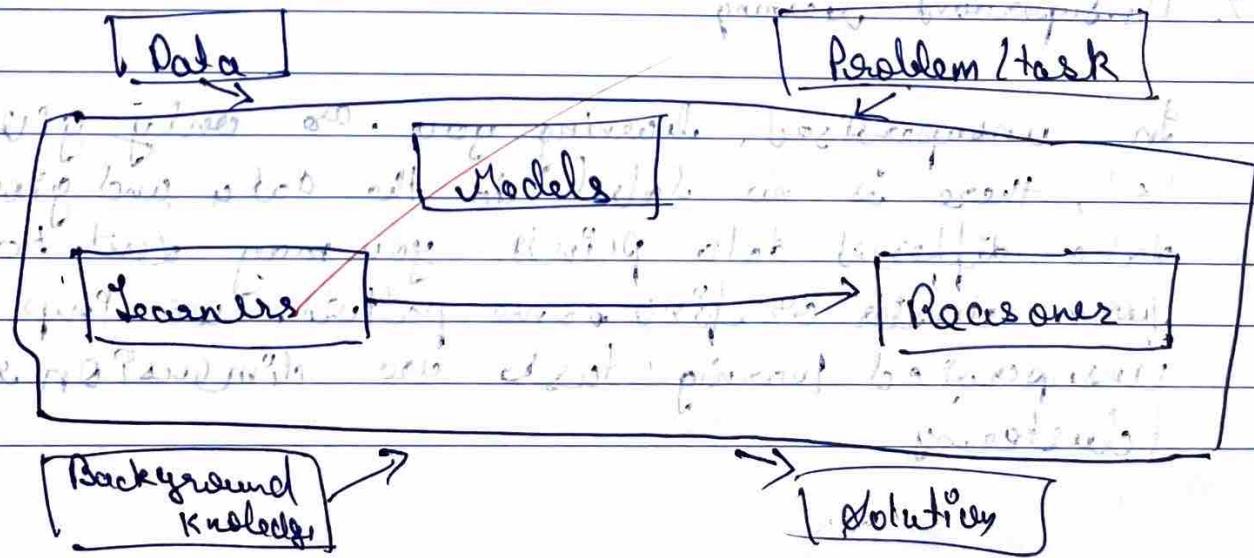
without

## Assignment 3

Q1) What is learning? Explain different types of learning with example.

⇒ Machine learning is a category of artificial intelligence. In machine learning computers have the ability to learn themselves. Explicit programming is not required. Machine learning focuses on the study and development of algorithms that can learn themselves, explicit programming is not required. Machine learning focuses on the study and development of algorithm that can learn from data and also make predictions on data.

Machine learning mainly focuses on the design and development of computer programs that can teach themselves to grow and change when exposed to new data. Using machine learning we can collect information from a dataset by asking the computer to make some sense from data. ML is turning data into information.

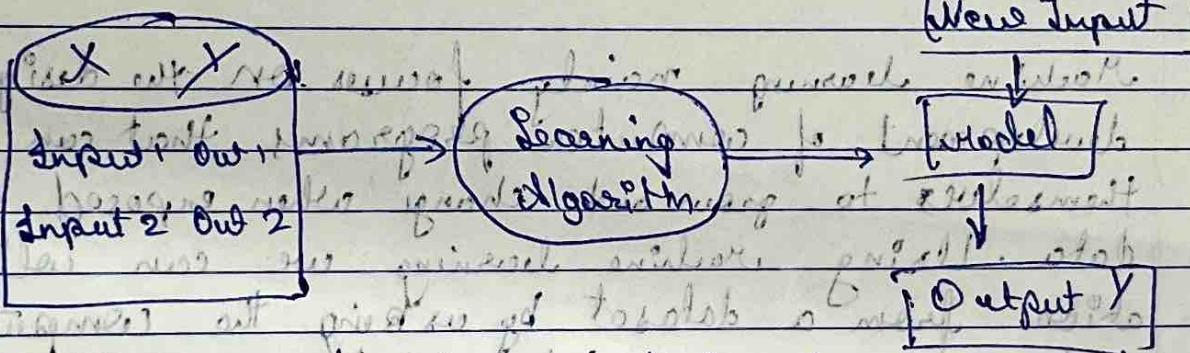


## types of Machine Learning

### 1. Supervised Learning

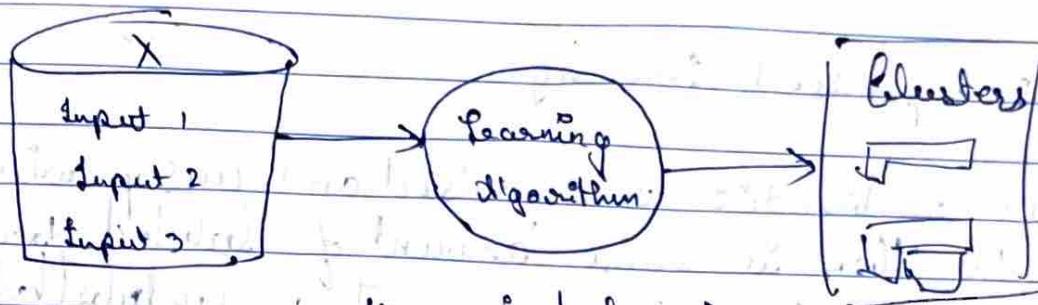
In this type of learning we use data which is comprises of input and corresponding output. In supervised learning training data is labelled with the correct answers e.g. 'spam' or 'ham'.

Two most important types of supervised learning are classification (where the output are discrete labels e.g. spam filtering) and regression (where the output are real values).



### 2. Unsupervised Learning

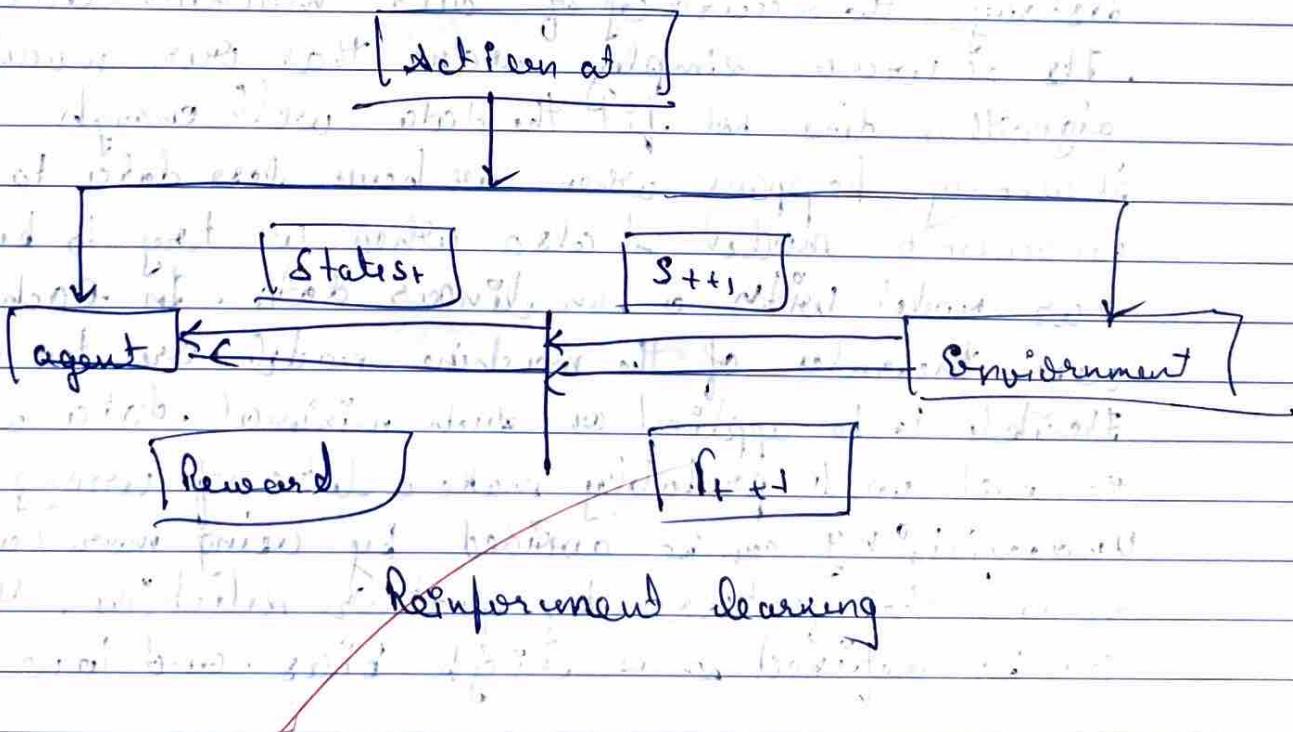
In unsupervised learning you are only given input  $x$ , there is no label to the data and given the data different data points, you may sent to form cluster or find some pattern. Two main unsupervised learning tasks are dimension reduction & clustering.



unsupervised learning

### 3. Reinforcement Learning

In reinforcement learning you have an agent who is acting in an environment and you want to find out what action the agent must take based on the reward or penalty that the agent gets. In this an agent, Example a robot or controller, acts to learn the optimal actions based on the outcomes of past actions.

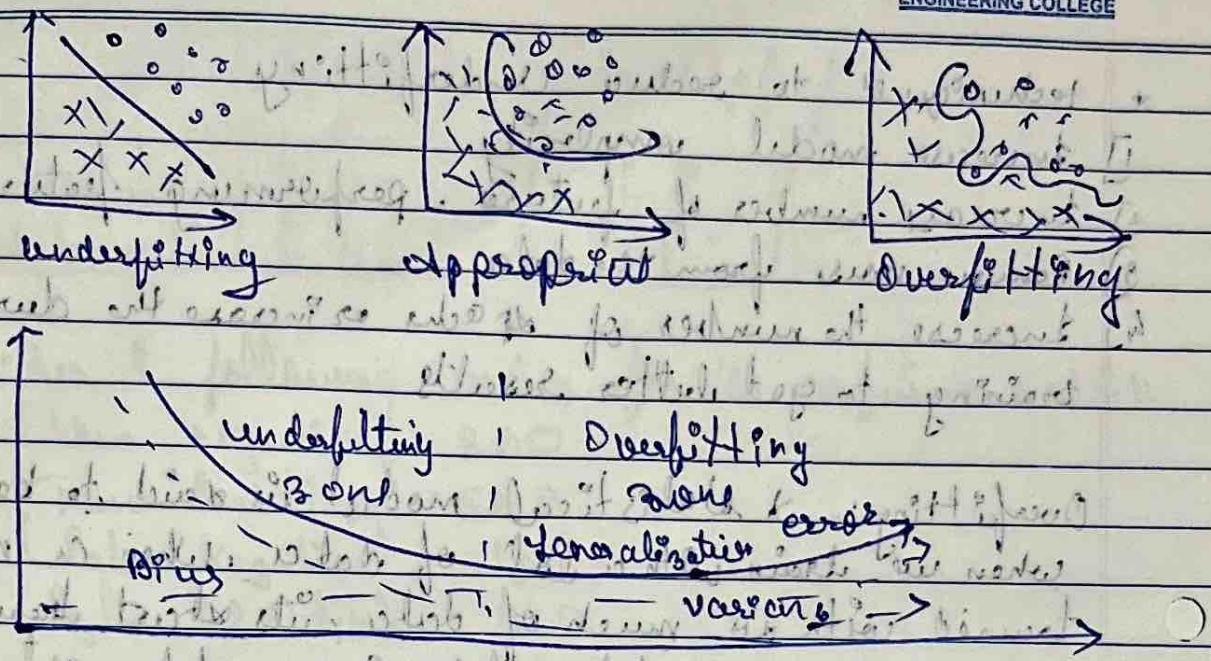


## q) Semi-Supervised Learning

It is a combination of supervised and unsupervised learning. In this there is some amount of labeled train data and also you have large amount of unlabeled data and you try to come up with some learning algorithm that converts or even when training data is not labeled by test classification image classification, where small amount of labeled

## q2) Define Overfitting &amp; Underfitting.

⇒ Underfitting A statistical model or a machine learning algorithm is said to have underfitting when it cannot capture the underlying trend of the data. Underfitting destroys the accuracy of our machine learning model. Its occurrence simply means that our model on the algorithm does not fit the data well enough. It usually happens when we have less data to build → an accurate model. It also when we try to build a linear model with non-linear data. In such such case the rules of the machine model are too easy & flexible to be applied on such minimal data and therefore the model will probably make a lot of wrong prediction. Underfitting can be avoided by using more data & also reducing the features by feature selection. Underfitting can be refined as high bias & low var.



a) Illustrate process of learning with the gradient descent for a univariate linear regression using a bell shaped error function. Explain how it step by step updates the model parameters on every iteration.

⇒ Gradient descent is an iterative optimization algorithm that tries to find the optimum value (minimum or maximum) of an objective function. It is one of the most used optimization techniques in machine learning projects for updating the parameters of a model in order to minimize a cost function.

The main aim of gradient descent is to find the best parameter of a model which gives the highest accuracy of training as well as testing datasets. In gradient descent, the gradient is a vector that points in the direction of the steepest increase of the function at a specific position. Moving in the opposite direction of the gradient allows the algorithm to gradually descend towards lower values of the function and eventually reaching to the minimum of the function.

- \* techniques to reduce Underfitting
  - 1) Increase model complexity.
  - 2) Increase number of features, performing feature engineering.
  - 3) Remove noise from the data.
  - 4) Increase the number of epochs or increase the duration of training to get better results.

Overfitting: A statistical model is said to be overfitted when we train with a lot of data, when a model gets trained with so much of data, it's start learning from the noise and inaccurate data entries in our data set. Then the model does not categorize the data correctly because of too many details and noise. The cause of overfitting are the non-parametric and non-linear methods because these types of machine learning algorithms have more freedom in building the model based on the data set and therefore they can really build unrealistic models.

A solution to avoid overfitting is using a linear algorithm if we have linear data or using the parameters like maximal depth if we are using decision trees.

Overfitting: It can be defined as fitting and allowing

any technique to reduce overfitting.

- 1) Increase training data.
- 2) Reduce model complexity (like early stopping).
- 3) Early stopping during the training phase.
- 4) Ridge Regularization.
- 5) Use dropout for neural network to tackle overfitting.

	Predicted	Predicted
Actual 0	T N	F P
Actual 1	F N	T P

Consider the following values for the confusion matrix

$$TN (\text{true negative}) = 300$$

$$TP (\text{true positive}) = 500$$

$$FN (\text{false negative}) = 150$$

$$FP (\text{false positive}) = 60$$

### a) Accuracy

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

accuracy is defined as the ratio of the number of correct predictions and the total number of predictions. It lies between (0, 1). Higher accuracy means a better model ( $TP + TN$ ). Higher accuracy means better model. ( $TP + TN$ ) must be high.

Accuracy is not a useful metric in case of an imbalanced datasets, e.g. 960 patients without cancer & 10 patient with cancer

$$\text{accuracy using above values} : \frac{(500+300)}{(500+30+150+60)} = 800 \cdot 800\%$$

### \* Precision

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

Precision is useful when we want to reduce the number of positives. Eg. a system that predicts if the email received is spam or not taking spam as a positive, then we don't want our system to predict non-spam emails.

Precision using above values

$$\text{Precision} = \frac{500}{500 + 50} = \frac{500}{550} = 90.90\%$$

d) Recall

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{500}{500 + 10} = \frac{500}{510} = 98.04\%$$

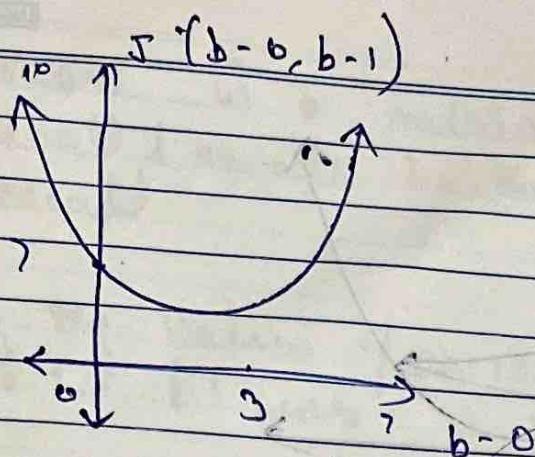
Recall is useful in case of cancer detection, where we want to minimize the number of false negatives for any practical use since we don't want to work person with cancer as safe.

Recall using above values

$$\text{Recall} = \frac{500}{500 + 10} = \frac{500}{510} = 98.04\%$$

e) F-1 score

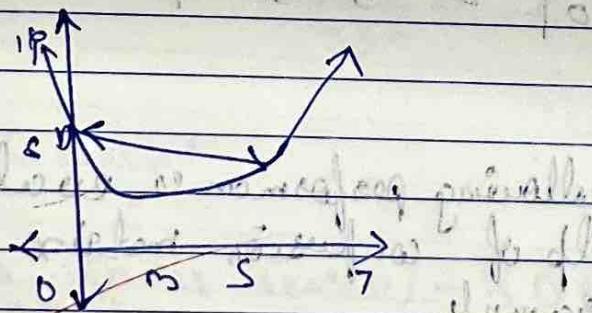
$$F1 = \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$



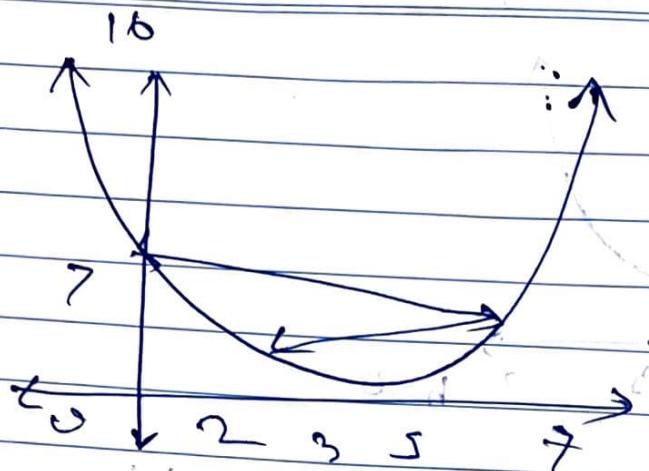
As we can see we have a simple parabola with a minimum at  $b = 0 = 3$ . The model will use this gradient descent algorithm to gradually tune the values of  $b = 0$  &  $b = 1$ .

To start this process, our model will initialize all parameters to 6 this means our initial cost is  $J$

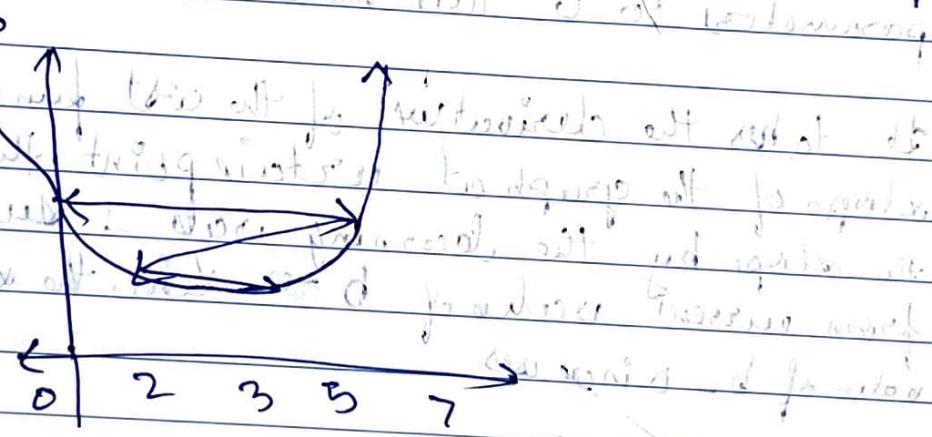
It takes the derivatives of the cost function & returns the slope of the graph at certain point, the model multiplies the slope by the learning rate & subtracts the product from current value of  $b = 0$ . Since the slope is negative, the value of  $b = 0$  increases.



The value of  $b = 0$  increases the algorithm will now take the derivative of cost function when  $b = 0$  is 5 & carry out the process again.



This process will continue until the derivative is returning a slope that is very close to 0, signifying that the model has reached a number of optimal values for the parameters being tested.



A) Explain the following performance evaluation parameters with the help of confusion matrix. Illustrate with appropriate example

⇒ A confusion matrix is a table that is often used to describe the performance of a classification model (or classifier) on a set of test data for which the true values are known.

F1 score is a metric that combines precision and recall & equals to the harmonic mean of precision and recall.

Using the values of precision & Recall (0.9040 & 0.7692)  $F1\text{ score} = \frac{2 \cdot 0.9040 \cdot 0.7692}{0.9040 + 0.7692} = 0.8333 = 83.33\%$ .

### c) Specificity

Specificity is defined as the ratio of True Negatives & True negatives + False Positives. We want its value of specificity to be high.

specificity =  $\frac{\text{True Negative}}{\text{True Negative} + \text{False Positive}}$

thus specificity using above values will be

$$\frac{300}{(300 + 50)} = \frac{300}{350} = 85.71\%$$

### d) AUC - ROC curve

AUC (Area Under Curve) - ROC (Receiver Operating Characteristic) curve is one of the most important evaluation metrics for checking any classification models performance.

It is plotted between FPR (X-axis) & TPR (Y-axis). If the value is less than 0.5 then the model is even worse than random.

queering model solution to the problem  
maximum moment at the strong & moderate base  
(lower base)

$$TPR = TP$$

top of model  $\frac{FP + FN}{L}$  for middle stiff part

$$FR = \frac{FP}{FP + TN}$$

$$\frac{FP + TN}{L}$$

bottom part to write it in terms of stiffness  
middle part, writing out first & middle part &  
first end at  $\frac{FP + TN}{L}$  for middle  
(A) ~~stiff~~  
 middle part, writing out  
 middle part, writing out

1. Middle section ends process of writing out

$$\times 15.28 - 0.02 : 62 + 0.081008$$

0.08

Hours 208 - 204 (9)

(Aug 2023) 208 - (Aug 2023) 204  
took away from for use in view (, 3 hours, 0.081008  
middle part, was middle part, written middle part  
 middle part, written middle part

(Aug 2023) 208 recorded bottom of it  
as middle part, the value of the (Aug 2023) 204 &  
middle part, middle part, written in, wrote out

## Assignment 2

d) Discuss different ensemble learning techniques

- Ensemble techniques combine individual models together, to improve the stability and predictive power of the model.
- Individual models, also known as base models can have a high bias, high variance or both, which is not good such models are known as weak models.
- When weak models are combined they can produce strong models having low bias and low variance.
  - A diverse set of models in comparison to single models can make better predictions. This diversification in ML is achieved by a technique called ensemble learning.
  - Ensemble methods → bagging, boosting

Some ensemble learning techniques are

- a) ~~Majority Voting~~: this method is generally used for classification problems. In this technique, multiple models are used to make prediction by each model are considered as a vote. The prediction we get from the majority of the models is used as a final prediction.
- b) ~~Averaging~~: Multiple predictions are made after each data point and the average of all prediction is the final prediction.
- c) ~~Weighted averaging~~: this is an extension of

of averaging where all models are assigned weight

Boosting :- It trains multiple models sequentially, one by one with each new model focusing on the mistakes made by the previous model, focusing on the mistakes made by the previous model. The idea is to correct the errors earlier models, the idea is to correct the errors earlier models & improve the overall prediction accuracy. The final model is usually a weighted average of all the models. It is usually a weighted average of all the models predictions. Boosting can significantly improve model performance and reduce both bias & variance. This is particularly useful for complex datasets.

- Q2) Explain the following  
i) Weak learners & strong learners  
→ A weak learner is a model that performs slightly better than random guessing. It is typically a simple model with limited predictive power, often making predictions that are only marginally better than random chance. In decision trees algorithms, a shallow tree with a few nodes and a simple rule-based classifier, can be considered a weak learner & strong learner is a model that performs significantly better than random guessing. It has a high predictive accuracy and it can capture complex patterns with depth. Eg:- Complex models like deep neural networks

Strong learners are often the goal ensemble methods. They are built by combining weaker models or using sophisticated algo to achieve high per dictive performance.

Metac learning :  
It refers to the process of learning how to learn. It involves developing models that can learn from various learning algorithms or datasets and adapt to new tasks more efficiently. The goal is to improve the learning process itself by using previous experience or data. e.g. Meta learning can involve learning the best hyperparameters for different algorithms, collecting appropriate models for various tasks or even adapting to new types of data with minimal training.

Random forest :  
It is an ensemble learning method specifically designed for classification & regression tasks. It builds a number of decision trees & combines their predictions to make a final decision.  
Working Bagging - Random forest uses the Bagging technique to train each decision tree on a different random subset of training data & upon running the training of each decision tree after regression tasks, the final prediction is the average of the prediction of all trees.

Q3) Compare Bagging, Boosting & Stacking

Bagging

Boosting

Stacking

- 1) Parallel training (independent models) → sequential training (combines correct errors) → ensemble of trees
- 2) Average or majority voting! → weighted average based on model's base model prediction

- 3) Reduces variance → Reduces both bias and variance of individual models

- 4) less prone to overfitting → prone to overfitting

- 5) less computational → can be computed sequentially

- 6) simpler to implement → More complex, many models understood sequentially

- 7) E.g.: Random forest vs e.g.: Adaboost → E.g.: Model stacking → multiple models → further decisions → trees, logistic regression

Q) Explain Adaboost algorithm

→ Adaboost is a boosting algorithm that also works on the principle of stage wise addition method where multiple weak learner models are used to get strong learner models. The value of alpha parameters in this case will be indirectly proportional to the errors of weak learners.

- Working

- 1) Initialization :- All training samples are given equal weight.

- 2) Training : The algorithm iteratively trains weak classifiers

- 3) Weighting : The influence of each weak learner is determined by its accuracy

- 4) Updating : - Weights of misclassified samples are increased for the next iterations

- 5) Final model : - The final prediction is a weighted vote of all weak learners

- This algorithm enhances performance by focusing on outliers & it is generally more robust and to overfitting than other models

Q 4  
introduction

- Advantages
- simple & effective
- no overfitting
- can be used with various types of weak learners
- Disadvantages
- sensitive to noisy data
- It is computationally expensive