(Q1) In logistic regression, what is the logistic function (sigmoid function) and how is it used to compute probabilities?

In logistic regression, the logistic function, also known as the sigmoid function, is a mathematical function that maps any real-valued number to a value between 0 and 1

It is defined by the formula:

$$\sigma(z) = 1 + e^{-z}1$$

Where:

$\sigma(z)$ is the output of the sigmoid function

$z$ is the input to the sigmoid function

The sigmoid function has the following properties:

1.It produces values in the range [0, 1]

2.It is symmetric around the origin (0.5 when input is 0)

3.It is differentiable, which is important for gradient-based optimization algorithms

Computing Probabilities with Logistic Regression:

f(y)=1 /1+e-y

y=b0+b1x

After substitution formula will be: f(y)=1/1+e-(b0+b1x)

(Q2) When constructing a decision tree, what criterion is commonly used to split nodes, and how is it calculated?

**Gini impurity**: Gini impurity measures the probability of incorrectly classifying a randomly chosen element in the dataset if it were randomly labeled according to the class distribution. It is calculated by summing the squared probabilities of each class being chosen.

Mathematically, the Gini impurity $IG$ for a node with class probabilities $1,2,...,p1,p2,...,pk$ is calculated as:

$$IG = 1 - \sum i=1 kpi2$$

**Entropy**: Entropy measures the average amount of information needed to predict the class of a randomly chosen data point. It is calculated using the Shannon entropy formula.

Mathematically, the entropy $IH$ for a node with class probabilities $1,2,...,p1,p2,...,pk$ is calculated as:

$$IH = -\sum i=1 kpi\log2(pi)$$

(Q3) Explain the concept of entropy and information gain in the context of decision tree construction?

Aim: Increasing homogeneous groups based on their features.

Entropy: It measures the impurity or uncertainty within a node If, Entropy is high purity is low

Formula: Entropy=-(sum(pilog2(pi)))

Information Gain: It measures the reduction in uncertainty by splitting a node on a particular feature.  If, IG is high purity is high

 Formula: IG=Entropy(s)-sum(s1/s) * entropy(s1)

 While node splitting make sure node entropy should be less and IG should be high for best decision☐making process.

## (Q4) How does the random forest algorithm utilize bagging and feature randomization to improve classification accuracy?

Random forests improve classification accuracy by two techniques: bagging and feature randomization.

Bagging (Bootstrap Aggregation):

☐ Bagging also called as Bootstrap Aggregation.

☐ It involves creating multiple decision trees.

☐ A bootstrap sample is a random sample with replacement (some data points may appear multiple times while others are omitted).

☐ Improve accuracy of classification by VOTING technique

☐ The diversity in training data leads to individual trees with different biases and strengths. (it reduces the variance of the model and becomes less prone to overfitting. Hence, significantly improve accuracy on unseen data/new data.)

 Feature Randomization:

☐ Random subset of features is considered for splitting.

☐ Different trees may use different features to arrive at the same decision (further adding diversity to the ensemble).

☐ Improve accuracy by preventing any single feature from dominating the decision-making process

## (Q5) What distance metric is typically used in k-nearest neighbors (KNN) classification, and how does it impact the algorithm's performance?

The most commonly used distance metric in k-nearest neighbors (KNN) classification is the Euclidean distance. However, other distance metrics such as Manhattan distance, Minkowski distance, and cosine similarity can also be used depending on the nature of the data and the problem at hand.

**Euclidean Distance:**

Euclidean distance is the straight-line distance between two points in Euclidean space. It is calculated as the square root of the sum of the squared differences between corresponding coordinates of the two points.

Mathematically, for two points $p$ and $q$ in $n$-dimensional space, the Euclidean distance $d$ is calculated as:

$$d(p,q)=\sqrt{\sum_{i=1}^{n}(p_i-q_i)^2}$$

**Manhattan Distance:**

Manhattan distance, also known as city block distance or $L1$ norm, is the sum of the absolute differences between corresponding coordinates of two points.

Mathematically, for two points $p$ and $q$ in $n$-dimensional space, the Manhattan distance $d$ is calculated as:

$$d(p,q)=\sum_{i=1}^{n}|p_i-q_i|$$

**Minkowski Distance:**

Minkowski distance is a generalization of Euclidean and Manhattan distances. It is defined as:

$$d(p,q)=\left(\sum_{i=1}^{n}|p_i-q_i|^p\right)^{1/p}$$

When $p=2$, it reduces to Euclidean distance. When $p=1$, it reduces to Manhattan distance

## (Q6) Describe the Naïve-Bayes assumption of feature independence and its implications for classification?

Implications of the Feature Independence Assumptions:

☐ Simplified Computation: efficient calculation of probabilities, makingNaive Bayes a computationally inexpensive algorithm.

☐ Robustness to Missing Values: Naive Bayes can handle missing values in individual features without significant performance degradation, as it treats them as independent.

☐ Scalability to High Dimensions: Work effectively with large numbers of features, making it suitablefor high-dimensional dataset.

## (Q7) In SVMs, what is the role of the kernel function, and what are some commonly used kernel functions?

Role of Kernel Function:

☐ Takes two input data points.

☐ Applies a mathematical operation on them.

☐ Outputs a value that measures their similarity in the higher-dimensional feature space.

☐ This similarity is then used by the SVM to determine the optimal decision boundary for

classification.

Commonly used kernel functions:

☐ Linear kernel: The simplest kernel, suitable for linearly separable data.

☐ Polynomial kernel: Creates higher-order polynomial features from the original data, useful forcapturing complex non-linear relationships.

☐ Radial Basis Function (RBF kernel): A popular choice, uses a Gaussian function to

measuresimilarity, making it flexible for different data distributions.

☐ Sigmoid kernel: Similar to RBF, but with a different activation function, suitable for specific cases

## (Q8) Discuss the bias-variance trade off in the context of model complexity and overfitting?

The bias-variance tradeoff is a fundamental concept in machine learning that describes the relationship between the bias of a model and its variance, and how they affect the model's performance.

### Bias:

Bias refers to the error introduced by approximating a real-world problem with a simplified model. A high bias model tends to oversimplify the data and may not capture the underlying patterns effectively.

Models with high bias have low complexity and are often too rigid, leading to underfitting. They may perform poorly on both training and test data

### Variance:

Variance refers to the variability of model predictions for different training datasets. A high variance model is sensitive to small fluctuations in the training data and may capture noise instead of the underlying patterns

Models with high variance have high complexity and are often too flexible, leading to overfitting. They may perform well on the training data but generalize poorly to unseen data

### Tradeoff:

The bias-variance tradeoff arises from the fact that increasing model complexity typically reduces bias but increases variance, and vice versa

The bias-variance tradeoff arises from the fact that increasing model complexity typically reduces bias but increases variance, and vice versa

## (Q9) How does TensorFlow facilitate the creation and training of neural networks?

TensorFlow offers several powerful features that make it a popular choice for creating and training neural networks:

1.Ease of Use:

• High-level API

• Automatic differentiation

2.Flexibility and Customization:

• Low-level API/High API on our need

• Wide range of supported models (CNN, RNN etc..,)

3.Scalability and Efficiency:

• Tensor Board

• Distributed training

• TensorFlow Lite

4.Community and Resources:

• Large community

• Open-source

## (Q10) Explain the concept of cross-validation and its importance in evaluating model performance?

Cross-validation is used to evaluate the performanceof a model and prevent overfitting.

Working:

1. Data Split: Divide your dataset into folds (typically 5 or 10).

2. Train-Test Split: For each fold:

o Use k-1 folds as the training set to train your model.

o Use the remaining 1 fold as the testing set to evaluate the model's performance.

3. Repeat: Repeat steps 1 & 2 for all folds.

4. Evaluation: Calculate a performance metric (e.g., accuracy, precision, F1-score) for each fold.

5. Average: Take the average of the performance metrics across all folds to get an overall estimate ofthe model's performance.

Importance of cross-validation:

• Prevents overfitting

• Provides a more realistic estimate

• Enables comparing models

(Q11) What techniques can be employed to handle overfitting in machine learning models?

Overfitting is a model memorizes the training data too well when a model becomes too focused on the specific details of the training data (leading to poor performance on unseen data).

Techniques employ to handle overfitting:

• Increase Data Size

• Data Cleaning and Preprocessing

• Choose simpler models like decision trees or linear regression after that gradually increase

the complexity of model

• Regularization techniques like L1/L2

• Dropout encourages the model to learn more robust features.


Q12) What is the purpose of regularization in machine learning, and how does it work?

Regularization aim is to prevent overfitting and improve the generalizability of the model.

It helps avoid overfitting by introducing penalties or constraints that discourage the model from becoming overly complex or fitting the training data too closely.

Regularization works by adding a penalty or complexity term to the complex model.

 Regularized Cost=Original Cost + $\lambda \times$ Regularization Term

o original cost function (which measures how well the model fits the training data)

o regularization term (which penalizes complex models)

o $\lambda$ is the regularization parameter, which controls the strength of the regularization.

L1 Regularization (Lasso):

• It adds the absolute values of the coefficients as a penalty term to the cost function.

• The regularization term is proportional to the sum of the absolute values of the model

parameters.

• It tends to produce sparse models, meaning it encourages some of the coefficients to become exactly zero, effectively performing feature selection.

L2 Regularization (Ridge):

• It adds the squared values of the coefficients as a penalty term to the cost function.

• The regularization term is proportional to the sum of the squared values of the model

parameters.

• It tends to penalize large coefficients and can mitigate multicollinearity (high correlation

between predictor variables) by distributing the weight across correlated features.

(Q13) Describe the role of hyper-parameters in machine learning models and how they are tuned for optimal performance?

Hyperparameters are like the knobs and dials of the model.

These are external settings that control the learning process and model complexity, ultimately impacting its performance. Unlike parameters, which are learned from the data during training,hyperparameters are set before training begins.

Learning rate: Controls the step size taken during gradient descent optimization in algorithms like linear regression and neural networks.

• Number of hidden layers and neurons in neural networks: Affects the model's capacity to learn complex non-linear relationships.

• Regularization parameters: (e.g., L1/L2 regularization strength) Control the model's complexity to prevent overfitting.

• Kernel function in Support Vector Machines: Determines the way data points are compared in the feature space

(Q14) What are precision and recall, and how do they differ from accuracy in classification evaluation?

**Precision**:

Precision measures the proportion of true positive predictions among all positive predictions made by the model.

It is calculated as the ratio of true positives (correctly predicted positive instances) to the sum of true positives and false positives (instances incorrectly predicted as positive).

High precision means that when the model predicts a positive class, it is likely to be correct.

$$Precision = TP / (TP + FP)$$

**Recall**:

Recall, also known as sensitivity or true positive rate, measures the proportion of true positive predictions among all actual positive instances in the data.

It is calculated as the ratio of true positives to the sum of true positives and false negatives (instances incorrectly predicted as negative but are actually positive).

High recall means that the model correctly identifies most of the positive instances in the data.

$$Recall = TP / (TP + FN)$$

(Q15) Explain the ROC curve and how it is used to visualize the performance of binary classifiers?

The Receiver Operating Characteristic (ROC) curve is a graphical representation that illustrates the performance of a binary classification model across different thresholds. It plots the true positive rate (TPR), also known as recall or sensitivity, against the false positive rate (FPR), which is the ratio of false positives to the total number of actual negatives.

**True Positive Rate (TPR)**:

TPR measures the proportion of true positive predictions among all actual positive instances. It is calculated as:

$$TPR = \frac{True\ Positives}{True\ Positives + False\ Negatives}$$

**False Positive Rate (FPR)**:

FPR measures the proportion of false positive predictions among all actual negative instances. It is calculated as:

$$FPR = \frac{False\ Positives}{False\ Positives + True\ Negatives}$$