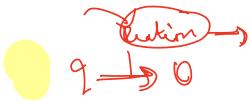


Q1.What Are the Different Types of Machine Learning?



Supervised Learning:

In supervised learning, the algorithm is trained on a **labeled dataset**, where each input example is paired with its corresponding output label. The goal is for the model to learn the **mapping** between **input features** and **output labels** so that it can make accurate predictions on new, **unseen data**.

Common Tasks:

Classification: Predicting categories or classes for new data points. For example, classifying emails as spam or not spam, or predicting whether a tumor is benign or malignant.

Regression: Predicting continuous numerical values. For instance, predicting house prices based on features like location, size, and number of bedrooms.

Unsupervised Learning:

Unsupervised learning involves training algorithms on **datasets without labeled responses**. Instead of learning from explicit feedback, the algorithm tries to identify patterns, relationships, or structures within the data.

Common Tasks:

Clustering: Grouping similar data points together based on some measure of similarity. For example, clustering customers into segments based on their purchasing behavior.

Dimensionality Reduction: Reducing the number of **input variables** or **features** while preserving the essential information. This is useful for visualizing high-dimensional data or speeding up subsequent computations.

Reinforcement Learning:

Reinforcement learning is about training agents to interact with an **environment** and learn from the consequences of their actions. The agent receives feedback in the form of rewards or penalties based on its actions, which guides its learning process.

Common Tasks:

Markov Decision Processes (MDPs): Reinforcement learning problems are often formulated as MDPs, where the agent takes actions in states and receives rewards accordingly.

Q2. What is Overfitting, and How Can You Avoid It?

The Overfitting is a situation that occurs when a model learns the training set too well, taking up random fluctuations in the training data as concepts. These impact the model's ability to generalize and don't apply to new data.

When a model is given the training data, it shows 100 percent accuracy—technically a slight loss.

But, when we use the test data, there may be an error and low efficiency. This condition is known as overfitting.

There are multiple ways of avoiding overfitting, such as:

- Penalty → Lasso Ridge
- Regularization: It involves a cost term for the features involved with the objective function
- Making a simple model :With lesser variables and parameters, the variance can be reduced
- Cross-validation methods like k-folds can also be used If some model parameters are likely to cause overfitting, techniques for regularization like LASSO can be used that penalize these parameters

Q3. What is 'training Set' and 'test Set' in a Machine Learning Model? How Much Data Will You Allocate for Your Training, Validation, and Test Sets?

Training Set:

The training set is a subset of the dataset used to train the machine learning model. It consists of input data paired with the corresponding correct output labels (in supervised learning). The model learns from this data by adjusting its parameters or weights through optimization algorithms (such as gradient descent) to minimize the error between its predictions and the actual labels.

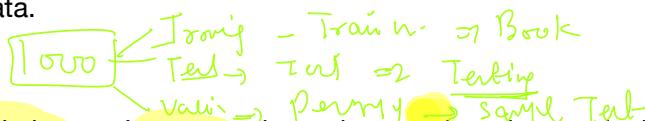
Validation Set:

The validation set is used to tune hyperparameters and evaluate the performance of the model during training. It helps in preventing overfitting by providing an unbiased evaluation of the model's performance on data that it hasn't seen during training. The validation set is also used for early stopping, where training is halted when the model's performance on the validation set stops improving or starts deteriorating.

Test Set:

The test set is used to evaluate the final performance of the trained model after it has been trained and validated. It provides an unbiased estimate of the model's performance on unseen data. The test set should ideally reflect the same distribution as the training and validation sets to ensure that the model generalizes well to new, unseen data.

Data Allocation:



The allocation of data to the training, validation, and test sets depends on various factors, including the size of the dataset, the complexity of the problem, and the available computational resources.

However, some common practices and guidelines are followed:

Training Set: Typically, the majority of the data is allocated to the training set, often around 60% to 80% of the total dataset. A larger training set allows the model to learn more effectively and generalize better to unseen data.

Validation Set: The validation set is usually smaller than the training set, typically around 10% to 20% of the total dataset. This subset is used for tuning hyperparameters and monitoring the model's performance during training.

Test Set: The test set is also smaller compared to the training set and is generally around 10% to 20% of the total dataset. It is kept completely separate from the training and validation sets until the final evaluation to ensure an unbiased estimate of the model's performance.

Q4. How Do You Handle Missing or Corrupted Data in a Dataset?

Remove Rows or Columns:

If the missing values are very few and removing them won't significantly affect the dataset's integrity, you can simply delete rows or columns containing missing values.

Mean/Median/Mode Imputation:

Replace missing values with the mean (average), median (middle value), or mode (most frequent value) of the respective feature. This method is straightforward and preserves the overall distribution of the data.

Forward Fill/Backward Fill:

Fill missing values with the value from the previous or next non-missing observation along the same feature. This is useful for time-series data where missing values are often consecutive.

Interpolation:

1 2 3

Use interpolation techniques such as linear interpolation or spline interpolation to estimate missing values based on neighboring data points. Interpolation is particularly useful for ordered or time-series data.

Imputation Models:

Train machine learning models (e.g., k-Nearest Neighbors, Random Forests, etc.) to predict missing values based on other features in the dataset. The model learns patterns from the complete data to impute missing values.

Flagging and Encoding:

Create an additional binary feature indicating whether a value is missing or not (flagging). Then, use encoding techniques such as mean imputation or model-based imputation for the missing values in that feature.

Q5. What Are the Three Stages of Building a Model in Machine Learning?

The three stages of building a machine learning model are:

Model Building Choose a suitable algorithm for the model and train it according to the requirement

Model Testing . Check the accuracy of the model through the test data

Applying the Model Make the required changes after testing and use the final model for real-time projects

Q6. What Are the Differences Between Machine Learning and Deep Learning?

Machine Learning:

Enables machines to take decisions on their own, based on past data

It needs only a small amount of data for training

Works well on the low-end system, so you don't need large machines

Most features need to be identified in advance and manually coded

The problem is divided into two parts and solved individually and then combined

Deep Learning :

Enables machines to take decisions with the help of artificial neural networks

It needs a large amount of training data

Needs high-end machines because it requires a lot of computing power

The machine learns the features from the data it is provided

The problem is solved in an end-to-end manner

Q7. What Are the Applications of Supervised Machine Learning in Modern Businesses?

Email Spam Detection

Here we train the model using historical data that consists of emails categorized as spam or not spam. This labeled information is fed as input to the model.

Healthcare Diagnosis

By providing images regarding a disease, a model can be trained to detect if a person is suffering from the disease or not.

Sentiment Analysis

This refers to the process of using algorithms to mine documents and determine whether they're positive, neutral, or negative in sentiment.

Fraud Detection

By training the model to identify suspicious patterns, we can detect instances of possible fraud.

Q8. What is Semi-supervised Machine Learning?

Supervised learning uses data that is completely labeled, whereas unsupervised learning uses no training data.

In the case of semi-supervised learning, the training data contains a small amount of labeled data and a large amount of unlabeled data.

Q9. What is the Difference Between Inductive Machine Learning and Deductive Machine Learning?

Aspect	Inductive Learning	Deductive Learning
Start with	Observations	A theory
Conclude with	A theory	Observations
Learning approach	Bottom-up	Top-down
Data need	Many specific examples	Established general rules
Generalization	Good at finding patterns	Excels at precise predictions
Handling Noise	Can be affected by noise	More robust with solid rules
Interpretation	Less clear due to data patterns	Easier to interpret, explicit rules
Thinking Skills	Creative, critical, inductive	Analytical, deductive
Applications	Real-life problem-solving	Abstract concept application
Reasoning	Goes from specifics to general	Goes from general to specific

Rules or logic to draw conclusion
 Rule-based system
 Expert systems
 Apply knowledge to solve problem.

Q10 .Compare K-means and KNN Algorithms.

Type of learning .	Supervised learning	Unsupervised learning
Task	Classification and regression	Clustering
Parameter	K, the number of nearest neighbors	K, the number of clusters
Input	Labeled data	Unlabeled data
Output .	Prediction or estimation of output variable based on k nearest neighbors .	Grouping of similar data points in k

Q11. What Is 'naive' in the Naive Bayes Classifier?

The classifier is called 'naive' because it makes assumptions that may or may not turn out to be correct.

The algorithm assumes that the presence of one feature of a class is not related to the presence of any other feature (absolute independence of features), given the class variable.

For instance, a fruit may be considered to be a cherry if it is red in color and round in shape, regardless of other features. This assumption may or may not be right (as an apple also matches the description).

Q12. How Will You Know Which Machine Learning Algorithm to Choose for Your Classification Problem?

While there is no fixed rule to choose an algorithm for a classification problem, you can follow these guidelines:

If accuracy is a concern, test different algorithms and cross-validate them

If the training dataset is small, use models that have low variance and high bias

If the training dataset is large, use models that have high variance and little bias

Q13 . When Will You Use Classification over Regression?

Classification is used when your target is categorical, while regression is used when your target variable is continuous. Both classification and regression belong to the category of supervised machine learning algorithms.

Examples of classification problems include:

Predicting yes or no

Estimating gender

Breed of an animal

Type of color

Examples of regression problems include:

Estimating sales and price of a product

Predicting the score of a team

Predicting the amount of rainfall

Q14 . Considering a Long List of Machine Learning Algorithms, given a Data Set, How Do You

Decide Which One to Use?

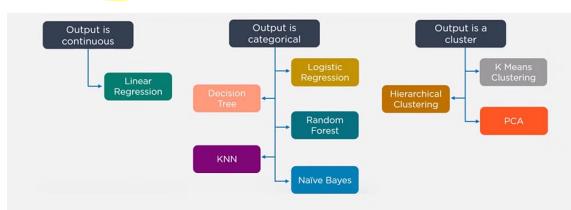
There is no master algorithm for all situations. Choosing an algorithm depends on the following questions:

How much data do you have, and is it continuous or categorical?

Is the problem related to classification, association, clustering, or regression?

Predefined variables (labeled), unlabeled, or mix?

What is the goal?



Q15.What is Bias and Variance in a Machine Learning Model?

Bias

Bias in a machine learning model occurs when the predicted values are further from the actual values. Low bias indicates a model where the prediction values are very close to the actual ones.

Underfitting: High bias can cause an algorithm to miss the relevant relations between features and target outputs.

Low Bias, Low Variance \rightarrow Underfitting

Variance

Variance refers to the amount the target model will change when trained with different training data.

For a good model, the variance should be minimized.

Overfitting: High variance can cause an algorithm to model the random noise in the training data rather than the intended outputs.

High variance, Low Bias = overfitting

Q16. What is the Trade-off Between Bias and Variance?

High Bias-Low Variance: Models with high bias and low variance tend to be simple and less flexible. They may not capture all the nuances in the data but are less affected by fluctuations in the training data. Examples include linear models or models with few parameters.

Low Bias-High Variance: Models with low bias and high variance tend to be more complex and flexible. They can capture intricate patterns in the data but are prone to overfitting, meaning they may perform well on the training data but generalize poorly to unseen data. Examples include decision trees with no pruning or deep neural networks.

The bias-variance trade-off in machine learning is about finding the right balance between simplicity and flexibility in a model. Too simple, and the model may miss important patterns (high bias); too complex, and it may overfit to noise in the data (high variance). The goal is to strike a balance that generalizes well to new data while capturing the underlying trends effectively.

Q17 . Achieving the balance between bias and variance involves several techniques:

Model Complexity: Start with a simple model and gradually increase complexity as needed. This

helps avoid overfitting initially.

Cross-Validation: Use techniques like k-fold cross-validation to evaluate model performance on different subsets of the data. This helps in understanding how well the model generalizes to unseen data.

Regularization: Introduce regularization techniques like L1 (Lasso) or L2 (Ridge) regularization to penalize overly complex models, discouraging them from fitting to noise.

Feature Selection: Select only the most relevant features to reduce the complexity of the model and prevent overfitting.

Ensemble Methods: Combine predictions from multiple models to reduce variance. Techniques like bagging (bootstrap aggregating) and boosting help in building robust models.

Hyperparameter Tuning: Tune model hyperparameters using techniques like grid search or random search to find the optimal settings that balance bias and variance.

Bias-Variance Decomposition: Understand the sources of error in the model by decomposing the overall error into bias and variance components. This helps in diagnosing model performance issues.

Q18. Define Precision and Recall.

Precision

Precision is the ratio of several events you can correctly recall to the total number of events you recall (mix of correct and wrong recalls).

$$\text{Precision} = (\text{True Positive}) / (\text{True Positive} + \text{False Positive})$$

Recall

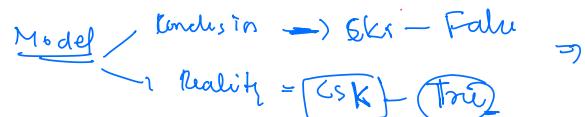
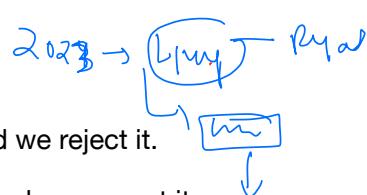
A recall is the ratio of the number of events you can recall to the number of total events.

$$\text{Recall} = (\text{True Positive}) / (\text{True Positive} + \text{False Negative})$$

Q19. What do you understand by Type I vs Type II error?

Type I Error: Type I error occurs when the null hypothesis is true and we reject it.

Type II Error: Type II error occurs when the null hypothesis is false and we accept it.



		reality			
		$H_0 = \text{True}$	$H_0 = \text{False}$		
Conclusion	H_0 is not rejected	OK	Type II error		
	H_0 is rejected	Type I error	OK		

		ACTUAL VALUES			
		Positive	Negative		
PREDICTED VALUES	Positive	TP	FP		
	Negative	FN	TN		

The predicted value is positive and its positive

Type I error : The predicted value is positive but it False

Type II error : The predicted value is negative but its positive

The predicted value is Negative and its Negative

Q21 .What is Cross-Validation?

Cross-validation is a technique used to evaluate the performance of a machine learning model by partitioning the dataset into subsets, training the model on a portion of the data, and then evaluating it on the remaining unseen data. This process is repeated multiple times, with different partitions of the data, and the performance metrics are averaged across all iterations.

The main steps of cross-validation are as follows:

Splitting the Data: The dataset is divided into k subsets of approximately equal size. One of these subsets is held out as the validation set, while the remaining $k-1$ subsets are used for training.

Training and Validation: The model is trained on the $k-1$ subsets and evaluated on the validation set. This process is repeated k times, each time using a different subset as the validation set.

Performance Evaluation: The performance metrics, such as accuracy, precision, recall, or mean squared error, are calculated for each iteration of the cross-validation process. These metrics are then averaged to obtain a more robust estimate of the model's performance.

Q22. What are the assumptions you need to take before starting with linear regression?

There are primarily 5 assumptions for a Linear Regression model:

Multivariate normality

No auto-correlation

Homoscedasticity

Linear relationship

No or little multicollinearity

linearly

Regularization ↗ Penalty

Q23. What is the difference between Lasso and Ridge regression?

The key difference is in how they assign penalties to the coefficients:

Ridge Regression:

Performs L2 regularization, i.e., adds penalty equivalent to the square of the magnitude of coefficients

Minimization objective = $\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^p \beta_j^2$ LS Obj + $\alpha * (\text{sum of square of coefficients})$

Lasso Regression:

Performs L1 regularization, i.e., adds penalty equivalent to the absolute value of the magnitude of coefficients

Minimization objective = $\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^p |\beta_j|$ LS Obj + $\alpha * (\text{sum of the absolute value of coefficients})$

Q24. Key Difference Between Ridge Regression and LASSO Regression

Feature Selection: Lasso can set coefficients to zero, effectively performing feature selection, while Ridge can only shrink coefficients close to zero.

Bias-Variance Tradeoff: Both methods introduce bias into estimates but reduce variance, potentially leading to better overall model predictions.

Regularization Technique: Ridge uses L2 regularization (squares of coefficients), and Lasso uses L1 regularization (absolute values of coefficients).

Predictive Performance: The choice between Ridge and Lasso depends on the data and the problem. Ridge tends to perform better for many significant predictors, while Lasso is more effective when only a few predictors are actually significant.

Q25. What do you understand by the F1 score? *Auc¹⁴*

$$\begin{aligned} F1 \text{ Score} &= \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}} \\ &= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \end{aligned}$$

Harmomic Mean Preciso & Recall

$$F_1 = \frac{2 \times P \times R}{P + R}$$

The F1 score is a metric commonly used to evaluate the performance of a classification model. It considers both the precision and recall of the model to provide a single numerical value that summarizes its effectiveness.

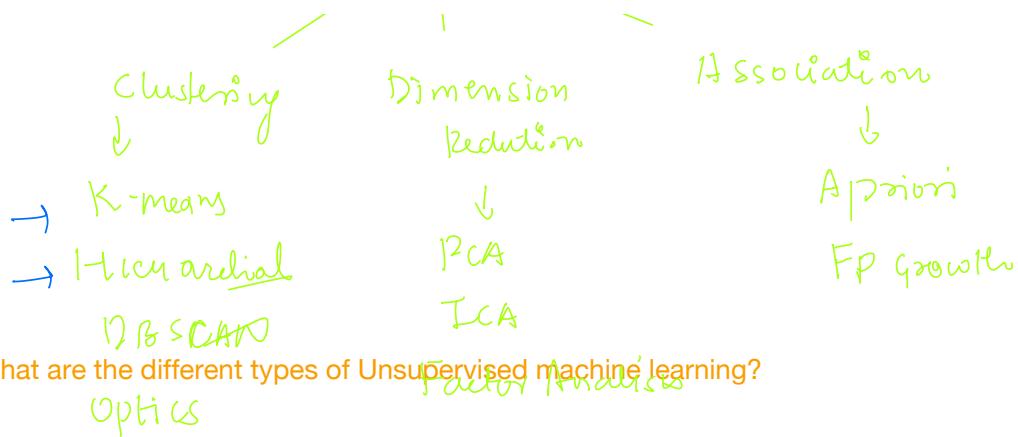
Precision: Precision measures the proportion of true positive predictions out of all positive predictions made by the model.

Recall (or Sensitivity): Recall measures the proportion of true positive predictions out of all actual positive instances in the dataset.

<u>Supervised</u>	<u>Unsupervised</u>	<u>Reinforcement</u>
<u>Target Value</u>	<u>Target Value</u>	<u>Objective Based</u>
Present	Not present	

Q26 .What are the different types of machine learning?





+

Q28. How do you know which machine learning algorithm should be used?

Regression: Predicting a Number (Linear Regression, Decision Trees, Random Forests, XGboost, Adaboost)

Classification: Predicting a Category (Logistic Regression, Decision Trees, Random Forests, XGboost, KNN, SVM)

Clustering: Grouping similar rows (K-Means, Hierarchical clustering, DBSCAN, OPTICS)

Dimension reduction: Reducing the number of variables in data (PCA, ICA, T-SNA, UMAP)

Association: Finding out which products sell together (Apriori, Eclat, FP-Growth)

Q29 . Can you explain the concept of overfitting in supervised learning?

Overfitting is a common problem in supervised machine learning where a model learns to capture the noise and random fluctuations in the training data instead of the underlying pattern or trend. This

leads to poor generalization performance, meaning the model performs well on the training data but fails to generalize to unseen data.

Issues:

Overfitting occurs when the model is too complex for the available data, capturing noise instead of patterns.

It happens when the model fits the training data too closely, including noise and random fluctuations.

High variance accompanies overfitting, indicating the model's predictions fluctuate widely with data changes.

Signs of overfitting include low training error but high validation error, poor performance on new data, and significant gaps between training and validation metrics.

Q30 . prevent overfitting ?

Regularization: Introducing penalties on the model parameters to prevent them from becoming too large.

Cross-validation: Evaluating the model's performance on multiple subsets of the data to ensure it generalizes well.

$$\text{Accuracy} = \frac{\text{No. of Correct predictions}}{\text{Total No. of Predictions}}$$

Feature selection: Selecting only the most relevant features to reduce the complexity of the model.

Using simpler models: Choosing a simpler model architecture that is less prone to overfitting.

Q31 . What evaluation metrics would you use to assess the performance of a classification model?

Accuracy: Accuracy measures the proportion of correctly classified instances out of all instances. It is calculated as the ratio of the number of correct predictions to the total number of predictions made .

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

Precision: Precision measures the proportion of true positive predictions out of all positive predictions made by the model. It is calculated as the ratio of true positives to the sum of true

positives and false positives.

$$\frac{TP}{TP + FN}$$

Recall (Sensitivity): Recall measures the proportion of true positive predictions out of all actual positive instances in the dataset. It is calculated as the ratio of true positives to the sum of true positives and false negatives.

F1 Score: The F1 score is the harmonic mean of precision and recall, providing a balanced measure of a model's performance. It is calculated as:

$$\frac{2 \times P \times R}{P + R}$$

ROC Curve and AUC: The Receiver Operating Characteristic (ROC) curve plots the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The Area Under the ROC Curve (AUC) provides a single scalar value representing the model's ability to discriminate between positive and negative instances.

Confusion Matrix: A confusion matrix provides a tabular summary of the model's predictions versus the actual class labels, showing true positives, true negatives, false positives, and false negatives.

Q33. What are some common techniques for feature selection in supervised learning?

Filter Methods: Evaluate feature relevance before model training based on statistical properties or correlation with the target variable. Techniques include:

Correlation Analysis

Chi-Square Test

Information Gain

10000 → Column

Wrapper Methods: Use model performance as a criterion for feature selection. Techniques include:

Forward Selection

Backward Elimination

Embedded Methods: Integrate feature selection into the model training process.

Lasso Regression

Tree-based Feature Importance

Dimensionality Reduction Techniques: Transform the feature space into a lower-dimensional space while preserving important information.

Principal Component Analysis (PCA)

Linear Discriminant Analysis (LDA)

Q34. How does regularization help prevent overfitting in supervised learning models?

Regularization and Overfitting:

- Regularization techniques introduce additional constraints or penalties on the model parameters during training to prevent overfitting.
- By penalizing large parameter values, regularization discourages the model from fitting the noise or random fluctuations in the training data.
- Regularization helps control the complexity of the model, ensuring that it generalizes well to unseen data by striking a balance between bias and variance.
- Common regularization techniques include L1 (Lasso) and L2 (Ridge) regularization, which add penalties to the absolute values and squared values of the model parameters, respectively.

Q35 . Can you describe the process of cross-validation and its importance in supervised learning?

Cross-Validation:

1. Cross-validation is a resampling technique used to evaluate the performance of a machine learning model on unseen data.
2. The process involves dividing the dataset into multiple subsets or folds, training the model on a portion of the data, and evaluating it on the remaining unseen data.
3. This process is repeated multiple times, with different partitions of the data, and the performance metrics are averaged across all iterations to provide a more robust estimate of the model's performance.

4. Common types of cross-validation include k-fold cross-validation, leave-one-out cross-validation, and stratified cross-validation, each with its own variations and applications.
5. Cross-validation is important in supervised learning because it helps assess how well the model generalizes to new, unseen data and provides insights into its stability and reliability.
6. It helps detect issues like overfitting or underfitting by evaluating the model's performance on multiple subsets of the data, enabling the selection of the best-performing model and hyperparameters.

Q35. What is the purpose of ensemble learning in supervised machine learning?

The purpose of ensemble learning in supervised machine learning is to improve predictive performance by combining the predictions of multiple individual models. Ensemble learning techniques leverage the diversity among these models to make more accurate and robust predictions compared to any single model alone.

Key purposes and benefits of ensemble learning include:

Improved Accuracy: Ensemble methods often achieve higher predictive accuracy compared to individual models by leveraging the wisdom of crowds. By combining multiple models, ensemble learning reduces the risk of selecting a suboptimal model and can better capture the underlying patterns in the data.

Robustness: Ensemble methods are more robust to noisy data and outliers because they aggregate predictions from multiple models. They tend to be less susceptible to overfitting, as errors made by individual models may cancel out when combined.

Reduction of Bias and Variance: Ensemble learning can help strike a balance between bias and variance, leading to more stable predictions. For example, combining models with high bias and low variance (e.g., decision trees) with models with low bias and high variance (e.g., neural networks) can result in an ensemble with lower overall bias and variance.

Model Interpretability: Ensemble methods can sometimes offer better interpretability than complex individual models. For example, in bagging methods like Random Forests, feature importance can be derived from aggregating feature importance scores across multiple trees.

Versatility: Ensemble learning is versatile and can be applied to various types of models and learning tasks. It can be used with classification, regression, and even unsupervised learning tasks.

Q35.What is the purpose of ensemble learning in supervised machine learning?

The purpose of ensemble learning in supervised machine learning is to improve predictive performance by combining the predictions of multiple individual models. Ensemble learning techniques leverage the diversity among these models to make more accurate and robust predictions compared to any single model alone.

Key purposes and benefits of ensemble learning include:

Improved Accuracy: Ensemble methods often achieve higher predictive accuracy compared to individual models by leveraging the wisdom of crowds. By combining multiple models, ensemble learning reduces the risk of selecting a suboptimal model and can better capture the underlying patterns in the data.

Robustness: Ensemble methods are more robust to noisy data and outliers because they aggregate predictions from multiple models. They tend to be less susceptible to overfitting, as errors made by individual models may cancel out when combined.

Reduction of Bias and Variance: Ensemble learning can help strike a balance between bias and variance, leading to more stable predictions. For example, combining models with high bias and low variance (e.g., decision trees) with models with low bias and high variance (e.g., neural networks) can result in an ensemble with lower overall bias and variance.

Model Interpretability: Ensemble methods can sometimes offer better interpretability than complex individual models. For example, in bagging methods like Random Forests, feature importance can be derived from aggregating feature importance scores across multiple trees.

Versatility: Ensemble learning is versatile and can be applied to various types of models and learning tasks. It can be used with classification, regression, and even unsupervised learning tasks.



Q36 .Common ensemble learning techniques include:



Bagging (Bootstrap Aggregating): Constructs multiple models independently and combines their predictions through averaging or voting. Examples include Random Forests.

Boosting: Builds models sequentially, where each subsequent model focuses on correcting the errors of the previous model. Examples include AdaBoost, Gradient Boosting Machines (GBM), and XGBoost.



Stacking (Stacked Generalization): Combines predictions from multiple models using a meta-model, which learns to weigh the predictions of base models.

Q37. How do you choose between different supervised learning algorithms for a given problem?How do you choose between different supervised learning algorithms for a given problem?

For Linear Relationships: Linear regression, Logistic regression, Linear Support Vector Machines (SVM).

For Non-linear Relationships: Decision Trees, Random Forests, Gradient Boosting Machines (GBM), Neural Networks.

For High-dimensional Data: Regularized regression (e.g., Lasso, Ridge), Support Vector Machines (SVM), Ensemble methods (e.g., Random Forests, Gradient Boosting).

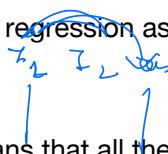
For Text Data: Naive Bayes, Support Vector Machines (SVM), Recurrent Neural Networks (RNNs), Transformer models (e.g., BERT).

Q38. What are the assumptions of Linear Regression?

Linearly →



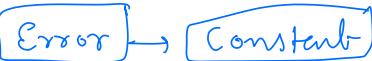
Linear Relationship between input and output – Linear relationship means if one increases, the other should also increase or vice-versa. Linear regression assumes that the input and output are linear dependents on each other.



No Multicollinearity – Multicollinearity means that all the input columns you have in data should not be highly related. For example, if we have X_1 , X_2 , and X_3 as input columns in data and if by changing X_1 there are changes observed in X_2 , then it is the scenario of multicollinearity. For example, if by

changing X1, we get output keeping that X2 and X3 are constant, but when X1 is correlated to X2, then on changing X1, X2 will also change, so we will not get proper output. That's why multicollinearity is a problem.

Normality of Residual – When we predict new data points and calculate error or residual (actual – predicted), then on plotting, it should be normally distributed across the mean. To know this, you can directly plot the KDE or QQ plot. First, you have to calculate the residual for every test point, and using the seaborn library; you can plot a distribution plot. The second way is to directly Plot a QQ plot where all points should be closer to the line.



Homoscedasticity – Home means same, and scedasticity means to spread or scatter. So this means having the same scatter. According to this assumption, when you plot residual, then the spread should be equal. If it is not equal, then it is known as Heteroscedasticity. To calculate this, you keep prediction on X-axis and residual on Y-axis. The scatter plot should be uniform.

No Autocorrelation of Errors – If you plot all residual errors, then there should not be a particular pattern.

Q38 .What are the key differences between logistic regression and linear regression ?

Problem Type:

Linear Regression: Linear regression is used for predicting continuous numeric outcomes. It models the relationship between the independent variables (features) and the continuous dependent variable (target) using a linear equation.

Logistic Regression: Logistic regression is used for predicting binary outcomes or probabilities. It models the probability that a given input belongs to a particular class using a logistic function, which restricts the output to be between 0 and 1.

Output:

Linear Regression: The output of linear regression is a continuous numeric value that represents the predicted outcome.

Logistic Regression: The output of logistic regression is a probability value between 0 and 1, which can be interpreted as the likelihood of the input belonging to a certain class. It is often used to classify inputs into one of two classes based on a threshold (e.g., 0.5).

Model Representation:

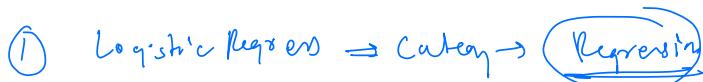
Linear Regression: In linear regression, the relationship between the input features and the target variable is represented by a straight line in a multidimensional space (hyperplane).

Logistic Regression: In logistic regression, the relationship between the input features and the log-odds of the target variable is represented by a sigmoid (logistic) function, which transforms the output to the range [0, 1].

Loss Function: → Error measure error function

Linear Regression: The loss function used in linear regression is typically the Mean Squared Error (MSE) or Mean Absolute Error (MAE), which measures the difference between the predicted and actual values.

Logistic Regression: The loss function used in logistic regression is the Binary Cross-Entropy (also known as Log Loss), which measures the difference between the predicted probabilities and the true binary labels.



Applications:

Linear Regression: Common applications of linear regression include predicting house prices, sales forecasting, and estimating the relationship between variables.

Logistic Regression: Common applications of logistic regression include binary classification tasks such as spam detection, disease diagnosis, and customer churn prediction.

Q40. How do support vector machines (SVMs) work in supervised learning? What is the kernel trick?

Basic Concept of SVM:

SVMs work by finding the optimal hyperplane that separates data points of different classes in a high-dimensional feature space.

The goal of SVMs is to maximize the margin, which is the distance between the hyperplane and the nearest data points (support vectors) of each class.

SVMs can handle both linearly separable and non-linearly separable data by mapping the input features into a higher-dimensional space where the classes become separable.

Linear Separable Data:

For linearly separable data, SVMs find the hyperplane that maximizes the margin while ensuring that all data points are correctly classified.

The hyperplane is determined by solving an optimization problem that involves minimizing the norm of the weight vector subject to the constraint that all data points are correctly classified.

Non-linearly Separable Data:

For non-linearly separable data, SVMs use a technique called the kernel trick to map the input features into a higher-dimensional space where the classes become separable.

The kernel trick allows SVMs to implicitly compute the dot product between the mapped feature vectors in the higher-dimensional space without explicitly transforming the data.

Kernel Trick:

The kernel trick is a mathematical technique that allows SVMs to operate in a high-dimensional feature space without explicitly computing the transformed feature vectors.

Instead of explicitly transforming the input features, the kernel function computes the dot product between the input feature vectors in the higher-dimensional space.

Q41., Bias and variance ?

Bias measures how well a model approximates the true relationship between features and target variable. High bias can cause underfitting.

Variance measures how much the model's predictions vary across different training datasets. High variance can cause overfitting.

Q42 , Regularisation :

Regularization is a technique used in machine learning and statistics to prevent overfitting and improve the generalization performance of models. It involves adding a penalty term to the model's loss function, which discourages the model from fitting the training data too closely or from becoming too complex.

The basic idea behind regularization is to impose constraints on the model parameters during training, preventing them from taking extreme or overly complex values. This helps to ensure that the

model captures the underlying patterns in the data without overfitting to noise or random fluctuations.

$$\rightarrow \text{Ridge} = (\lambda^2) \rightarrow \text{outlier} \rightarrow \text{Ridge}$$
$$\text{Lasso} = (\lambda) \rightarrow \text{outlier}$$

L1 Regularization (Lasso):

L1 regularization adds a penalty term proportional to the absolute values of the model parameters to the loss function.

This penalty encourages sparsity in the model by forcing some of the parameters to be exactly zero, effectively performing feature selection.

L1 regularization is useful when the dataset contains many irrelevant or redundant features.

L2 Regularization (Ridge):

L2 regularization adds a penalty term proportional to the squared values of the model parameters to the loss function.

This penalty encourages smaller parameter values, effectively shrinking the coefficients towards zero but rarely exactly to zero.

L2 regularization is useful for reducing the impact of multicollinearity and stabilizing the model's predictions.

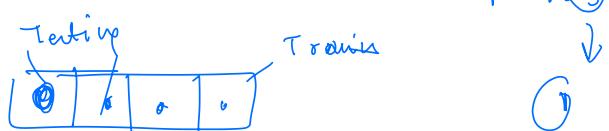
Q43 . cross validation :

Cross-validation is a resampling technique used in machine learning to evaluate the performance of a model on unseen data. It is commonly used to assess how well a predictive model generalizes to new data and to estimate its performance metrics.

The basic idea behind cross-validation is to partition the available dataset into multiple subsets or folds. The model is trained on a subset of the data, called the training set, and then evaluated on the remaining subset, called the validation set. This process is repeated multiple times, with different partitions of the data, and the performance metrics are averaged across all iterations to provide a more robust estimate of the model's performance.

There are several common types of cross-validation techniques:

K-Fold Cross-Validation:



The dataset is divided into K equally sized folds.

The model is trained K times, each time using $K-1$ folds for training and one fold for validation.

The performance metrics are averaged across all K iterations to obtain the final evaluation.

Leave-One-Out Cross-Validation (LOOCV):

Each data point is held out once as the validation set, and the model is trained on the remaining data.



This process is repeated for each data point in the dataset.

LOOCV is computationally expensive for large datasets but provides an unbiased estimate of the model's performance.

Stratified Cross-Validation:

Ensures that each fold has a similar distribution of classes or target variable values as the original dataset.

Particularly useful for imbalanced datasets where certain classes are underrepresented.

Q44. Hyperparameter ?

Hyperparameters are configuration settings that are external to the model and are not learned from the data during the training process. These parameters control the behavior of the learning algorithm and influence the performance and complexity of the model.

Definition:

Hyperparameters are parameters that are set prior to training and remain fixed throughout the training process.

They are distinct from model parameters, which are learned from the training data during the optimization process.

Examples:

Learning rate in gradient descent-based optimization algorithms.

Regularization parameter in L1 or L2 regularization techniques.

Number of hidden layers and units in a neural network.

Depth and number of trees in a decision tree-based model.

Kernel type and kernel parameters in support vector machines (SVMs).

Tuning:

Hyperparameter tuning, also known as hyperparameter optimization, involves selecting the optimal values for hyperparameters to improve the performance of the model.

This process typically involves conducting experiments with different combinations of hyperparameters and evaluating the model's performance using techniques like cross-validation.

Automated methods such as grid search, random search, and Bayesian optimization are commonly used for hyperparameter tuning.

Impact on Model:

Hyperparameters can significantly impact the performance, complexity, and generalization ability of the model.

Proper selection of hyperparameters is crucial for achieving the desired balance between bias and variance and for building models that generalize well to unseen data.

Suboptimal hyperparameters can lead to issues such as underfitting, overfitting, or poor model performance

Q44 . Accuracy in regression model :

Actual
predict → Error

Mean Absolute Error (MAE):

MAE measures the average absolute difference between the predicted and actual values.

It provides a straightforward interpretation of the average prediction error.

Mean Squared Error (MSE):

MSE measures the average squared difference between the predicted and actual values.

It penalizes larger errors more heavily than smaller errors and is widely used in optimization algorithms.

Root Mean Squared Error (RMSE):

RMSE is the square root of the MSE and provides a measure of the average magnitude of the errors.

It is in the same unit as the target variable and is easier to interpret than MSE.

Coefficient of Determination (R-squared):

R-squared measures the proportion of the variance in the dependent variable that is explained by the independent variables.

It ranges from 0 to 1, with higher values indicating a better fit of the model to the data.

Mean Absolute Percentage Error (MAPE):

MAPE measures the average percentage difference between the predicted and actual values relative to the actual values.

It provides a relative measure of accuracy and is useful for comparing models across different datasets.

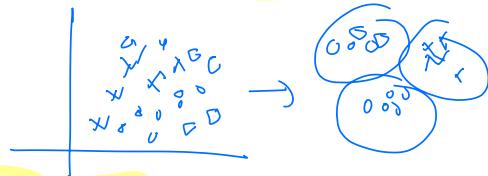
Adjusted R-squared:

Adjusted R-squared is a modified version of R-squared that penalizes the inclusion of unnecessary variables in the model.

It accounts for the number of predictors in the model and is more appropriate for comparing models with different numbers of variables.

Q45. Could you explain the K-means clustering algorithm and how it works? Explain this in same manner ?

Explanation of K-means Clustering:



K-means begins by randomly initializing K cluster centroids. These centroids represent the centers of the clusters. The number of clusters (K) is specified by the user beforehand.

Assignment Step:

Each data point in the dataset is assigned to the nearest cluster centroid based on some distance metric, typically Euclidean distance.

This step effectively groups data points into clusters by minimizing the distance between each point and its assigned centroid.

Update Step:

After all data points have been assigned to clusters, the centroids are recomputed based on the mean of all data points assigned to each cluster.

The centroids are moved to the average position of the points within their respective clusters.

This step iteratively updates the centroids to better represent the center of each cluster.

Convergence:

Steps 2 and 3 are repeated iteratively until either the centroids no longer change significantly or a predefined number of iterations is reached.

At convergence, the algorithm has effectively minimized the within-cluster sum of squares, meaning that the data points within each cluster are as close to the centroid as possible.

Output:

The final output of the K-means algorithm is a set of K cluster centroids and the cluster assignments for each data point.

Q46. What are the key challenges in determining the optimal number of clusters in K-means?

Elbow Method Ambiguity: It involves plotting the within-cluster sum of squares (WCSS) against the number of clusters and selecting the point where the rate of decrease in WCSS slows down.

However, this point may not always be clearly defined, making it challenging to determine the optimal Silhouette Analysis Interpretation: Silhouette analysis measures how similar an object is to its own cluster compared to other clusters. While silhouette scores provide insight into cluster quality, interpreting them can be subjective, and high average silhouette scores may not always lead to a clear choice of K.

Subjectivity and Context Dependency: Determining the optimal K can be subjective and context-dependent, as different stakeholders may have different perspectives on what constitutes meaningful clusters. Moreover, the optimal K may vary based on the specific problem domain and objectives.

Overfitting vs. Underfitting: Selecting too few clusters (underfitting) may oversimplify the data, while selecting too many clusters (overfitting) may lead to spurious or insignificant clusters. Balancing between underfitting and overfitting requires careful consideration and validation techniques.

High-Dimensional Data: In high-dimensional spaces, the distance between data points becomes less meaningful, affecting clustering performance. Determining an optimal K becomes more challenging due to the curse of dimensionality and increased data complexity.

Q47. K means vs Hierarchical clustering ?

Hierarchical Clustering:

Creates a hierarchy of clusters (dendrogram).

Does not require specifying the number of clusters beforehand.

Does not use centroids; builds hierarchy based on distance metrics.

No initialization required.

Computational complexity: $O(n^2 \log n)$ for agglomerative, $O(n^3)$ for divisive.

Handles clusters of arbitrary shapes and sizes.

K-means Clustering:

Partitions data into a predetermined number of clusters (K).

Requires specifying K beforehand.

Uses centroids to represent clusters; iteratively updates centroids.

Requires initialization of centroids.

Computational complexity: $O(n * K * I * d)$.

Assumes clusters are spherical and isotropic; may struggle with non-convex shapes.



Q48 . Can you explain the concept of dimensionality reduction and its importance in unsupervised learning?

Dimensionality reduction is the process of reducing the number of features or variables in a dataset while preserving its essential information. It is crucial in unsupervised learning because:

Efficiency: High-dimensional data can be computationally intensive and may suffer from the curse of dimensionality, leading to increased complexity and reduced performance of algorithms.

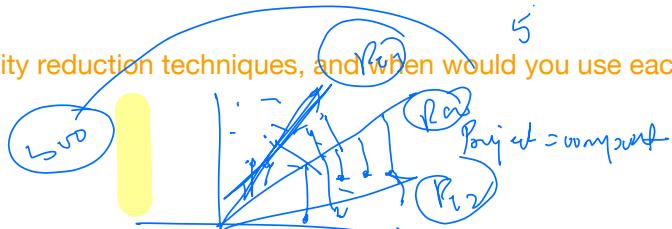
Visualization: Dimensionality reduction techniques like PCA and t-SNE help visualize complex datasets in lower-dimensional space, enabling better understanding and interpretation of data patterns.

Noise Reduction: By capturing the most relevant features and removing redundant or noisy ones, dimensionality reduction can improve the signal-to-noise ratio in the data, leading to better clustering or classification results.

Interpretability: Reduced-dimensional representations are often more interpretable and easier to analyze, facilitating insight discovery and decision-making processes.

Q48. What are some popular dimensionality reduction techniques, and when would you use each one?

Principal Component Analysis (PCA):



Use: PCA is widely used for linear dimensionality reduction. It identifies the principal components that capture the maximum variance in the data and projects the data onto these components.

When to use: PCA is suitable when the data has linear relationships between variables and when reducing dimensionality without losing much information is desired.

t-Distributed Stochastic Neighbor Embedding (t-SNE):

Use: t-SNE is effective for visualizing high-dimensional data in low-dimensional space (usually 2D or 3D) while preserving local similarities between data points.

When to use: t-SNE is suitable for visual exploration of data and clustering analysis, especially when the underlying data structure is non-linear.

Autoencoders:

Use: Autoencoders are neural network architectures that learn efficient representations of data by encoding it into a lower-dimensional latent space and then decoding it back to the original space.

When to use: Autoencoders are useful for nonlinear dimensionality reduction and feature learning, especially when the data has complex patterns and relationships.

Linear Discriminant Analysis (LDA):

Use: LDA is a supervised dimensionality reduction technique that maximizes the separation between classes while minimizing the variance within each class.

When to use: LDA is suitable when class information is available and when the goal is to reduce dimensionality while preserving class discrimination.

Kernel PCA:

Use: Kernel PCA extends PCA to nonlinear dimensionality reduction by using kernel methods to implicitly map data into high-dimensional feature spaces.

When to use: Kernel PCA is useful when the data has nonlinear relationships and when PCA may not capture the underlying structure effectively.

→ 1) Regression
→ 2) Clustering - Unsupervised
→ 3) Classification

Q49. How can you evaluate the performance of an unsupervised learning algorithm?

Silhouette Score: Silhouette analysis measures how similar an object is to its own cluster compared to other clusters. A higher silhouette score indicates better-defined clusters.

Davies-Bouldin Index: This index quantifies the average similarity between clusters, where lower values indicate better separation between clusters.

Calinski-Harabasz Index (Variance Ratio Criterion): This index measures the ratio of between-cluster dispersion to within-cluster dispersion. Higher values suggest better-defined clusters.

Internal Validation Metrics: Various internal metrics, such as inertia in K-means clustering or intra-cluster distances, can be used to assess the compactness and separation of clusters within the dataset.

External Validation Metrics: If some ground truth is available, external metrics like Adjusted Rand Index or Fowlkes-Mallows Index can be used to compare clustering results against known labels.

However, these are not always applicable in unsupervised settings.

Q50 .curse of dimensionality ?

The "curse of dimensionality" refers to the phenomena where the performance of certain algorithms degrades as the number of features or dimensions in the dataset increases. This concept is particularly relevant in high-dimensional spaces, where data points become increasingly sparse and distant from each other.

implications :

Empty Space

Data Sparsity: As the number of dimensions increases, the amount of available data per dimension decreases exponentially. This sparsity makes it challenging for unsupervised learning algorithms to identify meaningful patterns or clusters within the data.

Computational Complexity: Many unsupervised learning algorithms, such as clustering or dimensionality reduction techniques, rely on computing distances or similarities between data points. In high-dimensional spaces, the computational cost of these calculations increases significantly,

making algorithms computationally expensive and less scalable.

Overfitting: In high-dimensional spaces, the risk of overfitting also increases. Unsupervised learning algorithms may find spurious patterns or noise in the data, leading to overfitting and poor generalization performance on unseen data.

Difficulty in Visualization: Visualizing high-dimensional data becomes increasingly challenging as the number of dimensions increases.

Q51 . K means vs KNN ?

Purpose:

KNN: KNN is a supervised learning algorithm used for classification and regression tasks. It predicts the class or value of a new data point based on the majority class or average value of its nearest neighbors in the feature space.

K-means: K-means is an unsupervised learning algorithm used for clustering tasks. It partitions the data into a predefined number of clusters (K) based on the proximity of data points to cluster centroids.

Supervised vs. Unsupervised:

KNN: KNN is a supervised learning algorithm because it relies on labeled data to make predictions. It requires access to the ground truth labels during training.

K-means: K-means is an unsupervised learning algorithm because it does not require labeled data. It operates solely on the input features to identify clusters in the data.

Prediction vs. Clustering:

KNN: KNN predicts the class or value of a new data point by finding the K nearest neighbors in the feature space and using them to determine the outcome.

K-means: K-means clusters data points into K clusters based on their proximity to cluster centroids. It does not make predictions for new data points but rather assigns them to existing clusters based on their similarity to cluster centroids.

Number of Parameters:

KNN: KNN has hyperparameters such as the number of neighbors (K) and the choice of distance

metric.

K-means: K-means has hyperparameters such as the number of clusters (K) and the initialization method for cluster centroids.

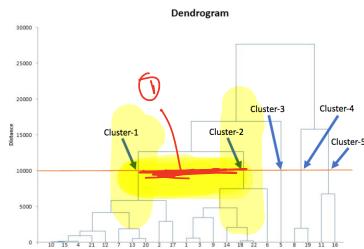
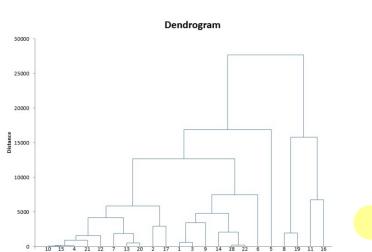
Algorithm Complexity:

KNN: KNN has a simple algorithmic complexity during inference, but it can be computationally expensive for large datasets since it requires calculating distances to all training samples.

K-means: K-means has a higher algorithmic complexity during training, as it involves iteratively updating cluster centroids until convergence. However, once trained, assigning new data points to clusters is computationally efficient.

Q52.: What are the types of hierarchical clustering?

Agglomerative Hierarchical Clustering: In agglomerative hierarchical clustering, each data point initially forms its own cluster, and then pairs of clusters are successively merged based on their proximity. This process continues until all data points belong to a single cluster. Agglomerative clustering can be visualized as a bottom-up approach, where smaller clusters are progressively merged into larger ones until a stopping criterion is met.



We cut the dendrogram at the point where the horizontal line covers the maximum distance between 2 vertical lines.

Divisive Hierarchical Clustering: In divisive hierarchical clustering, all data points initially belong to a single cluster, and then the algorithm recursively divides the cluster into smaller clusters based on some criteria, such as maximizing inter-cluster dissimilarity. This process continues until each data point forms its own cluster. Divisive clustering can be visualized as a top-down approach, where a single cluster is repeatedly split into smaller clusters until a stopping criterion is met.

Q53. how to decide the number of cluster is enough ?

Elbow Method: Plot the within-cluster sum of squares (WCSS) against the number of clusters (K). The "elbow" point on the plot, where the rate of decrease in WCSS slows down, can indicate the optimal number of clusters. This method suggests choosing the number of clusters at the point where adding more clusters does not significantly reduce the WCSS.

Silhouette Score: Compute the silhouette score for different values of K. The silhouette score measures the cohesion within clusters and the separation between clusters. A higher silhouette score suggests better-defined clusters. Select the value of K with the highest average silhouette score.

Q54. Steps for k means algorithm ?

Please look at the K-Means algorithm steps below:

1. If the value of K=2, then we need to find 2 groups in data
2. Assign group-1 or group-2 label to each observation in the data randomly
3. Find the “center” of all group-1 points and all group-2 points. These centers are known as “centroids”
4. Find the distances of each point from both the “centroids”
5. Assign each point to the nearest centroid. Hence, the corresponding group.
6. Repeat steps 3-5 till there are no further changes in the positions of the centroids OR the maximum number of iterations are completed.
7. In the end, each observation in data is assigned to one of the groups.

Q55. Linkages in hierarchical clustering ?

In hierarchical clustering, linkages determine how distances between clusters are calculated when merging or forming new clusters.

Single Linkage: Also known as nearest-neighbor linkage, it measures the distance between the closest points of two clusters when merging. It tends to create long, elongated clusters.

Complete Linkage: Also known as farthest-neighbor linkage, it measures the distance between the farthest points of two clusters when merging. It tends to create compact, spherical clusters.

Average Linkage: It calculates the average distance between all pairs of points from two clusters when merging. It balances between single and complete linkage, often resulting in balanced clusters.

Centroid Linkage: It calculates the distance between the centroids of two clusters when merging. It is less affected by outliers but may not capture the cluster structure accurately.

Q55. DBSCAN clustering :

DBSCAN stands for Density Based Spatial Clustering of Applications with Noise.

Terminologies :

1. **ϵ , epsilon (eps):** is known as the maximum allowed distance from one point to the other point, for both of them to be considered in one group/cluster
2. **MinPts:** is the minimum number of points which should be present close to each other at a distance of epsilon (ϵ) so that they all can form a group/cluster
3. **Core Point:** That point which has at least MinPts number of points near to it, within the distance of ϵ (eps)
4. **Border Point/Non-Core Point:** That point in data which has less than the minimum number of points(MinPts) within its reach (a distance of eps)
5. **Noise:** That point which has no point near to it within a distance of eps

DBSCAN steps :

1. Start with one point randomly, find out if it is a core point by checking the minimum number of points near to it by a distance of eps

2. If it is a core point, make it a cluster and move to the next unvisited point to repeat step-1
3. If the number of points within eps distance is fewer than MinPts then mark it as non-core point
4. If the number of points within eps distance is Zero, then mark that point as Noise.
5. Combine all those clusters together whose points are within eps distance. Also known as density connection or connected components. This starts a chain reaction. If cluster-1 and cluster-2 are connected and cluster-2 and cluster-3 are connected then cluster-1 and cluster-3 are also connected. Hence all these are combined to make one single cluster.

Q56. Factor analysis vs PCA ?

Factor Analysis tries to find the hidden groups of variables, which is like a common driving factor for the type of values in all those variables, e.g. customer survey, if you are unhappy with the taste of the coffee then all those questions about coffee taste in the survey will have low scores! Hence Bad Taste will be a Factor and questions like Coffee sweetness rating, Coffee bitterness rating, Coffee Freshness rating will represent the individual variables that will have low ratings.

Factor Analysis will form equations like below:

$$\text{Coffee sweetness rating} = \beta_1 * (\text{Bad Taste}) + C_1$$

$$\text{Coffee bitterness rating} = \beta_2 * (\text{Bad Taste}) + C_2$$

$$\text{Coffee Freshness rating} = \beta_3 * (\text{Bad Taste}) + C_3$$

PCA tries to find fewer new variables called Principal Components which are linear combinations of the variables in the data, hence these fewer variables “represent” all the variables, while reducing the dimensions. Here the goal is to create a new dataset that has the same number of rows but, lesser number of columns(Principal Components) which explains the variance of all the original variables in the data.

PCA will form equations like below:

$$PC_1 = \alpha_1 * (\text{Coffee sweetness rating}) + \alpha_2(\text{Coffee bitterness rating}) + \alpha_3(\text{Coffee Freshness rating})$$

$$PC_2 = \beta_1 * (\text{Coffee sweetness rating}) + \beta_2(\text{Coffee bitterness rating}) + \beta_3(\text{Coffee Freshness rating})$$

Q56 . Linear Regression vs Logistic regression?

Purpose:

Linear Regression: Linear regression is used for predicting continuous numerical outcomes. It models the relationship between one or more independent variables (predictors) and a continuous dependent variable (outcome) by fitting a linear equation to the observed data.

Logistic Regression: Logistic regression is used for predicting categorical outcomes, specifically binary outcomes (two classes). It models the probability of the presence of a certain outcome or event based on one or more independent variables, using a logistic (sigmoid) function to map predictions to the range [0, 1].

Output:

Linear Regression: The output of linear regression is a continuous numerical value. It predicts the value of the dependent variable based on the values of the independent variables.

Logistic Regression: The output of logistic regression is a probability score between 0 and 1. It represents the likelihood of belonging to a particular class or category.

Model Representation:

Linear Regression: The relationship between the independent variables and the dependent variable is represented by a straight line or hyperplane in the feature space.

Logistic Regression: The relationship between the independent variables and the log-odds of the dependent variable is represented by a sigmoid function, which produces an S-shaped curve.

Loss Function:

Linear Regression: Linear regression typically uses the mean squared error (MSE) or the sum of squared residuals as the loss function to minimize the difference between predicted and actual values.

Logistic Regression: Logistic regression uses the logistic loss (or log loss) function, also known as the cross-entropy loss, to minimize the difference between predicted probabilities and actual class labels.

Interpretation:

Linear Regression: The coefficients (slope and intercept) in linear regression represent the change in the dependent variable for a one-unit change in the independent variable.

Logistic Regression: The coefficients in logistic regression represent the log-odds ratio of the outcome variable for a one-unit change in the independent variable. They are interpreted in terms of odds ratios.

Logistic Regression is used for predicting a category, specially the Binary categories(Yes/No , 0/1).

When there are only two outcomes in Target Variable it is known as Binomial Logistic Regression.

If there are more than two outcomes in Target Variable it is known as Multinomial Logistic Regression.

Q58 .How a decision tree makes predictions:

A decision tree makes predictions by recursively splitting the dataset into subsets based on the values of input features. At each node, it selects the feature that best separates the data into purest subsets (homogeneous with respect to the target variable). This process continues until a stopping criterion is met, such as reaching a maximum depth or having a minimum number of samples in each leaf node. When a new data point is presented, it traverses the tree from the root node to a leaf node, where it predicts the majority class or average value of the training samples in that node.

Q59 .Advantages of decision trees:

1. Easy to interpret and visualize.
2. Can handle both numerical and categorical data.
3. Require little data preprocessing (e.g., feature scaling).
4. Non-parametric and robust to outliers.
5. Can capture nonlinear relationships between features and the target variable.

Q60 .Common criteria for splitting nodes in decision trees:

1. Gini impurity (Gini index).
2. Entropy.

3. Information gain (Gain ratio).

Q61. How pruning prevents decision trees from overfitting:

Pruning involves removing branches (subtrees) from the tree that do not provide significant improvements in predictive accuracy. It prevents decision trees from overfitting by reducing the complexity of the tree, thus promoting better generalization to unseen data.

Q62 .Can decision trees handle missing values in the dataset?:

Yes, decision trees can handle missing values in the dataset by using surrogate splits. Surrogate splits are alternative splits used when missing values are encountered during prediction. The decision tree algorithm automatically determines the best surrogate splits based on available data to make predictions for instances with missing values.

Q63 .What is a random forest and how does it work?:

A random forest is an ensemble learning method that builds multiple decision trees during training and outputs the mode (classification) or average prediction (regression) of the individual trees. Each tree in the random forest is trained on a bootstrap sample of the original dataset, and at each node, a random subset of features is considered for splitting. This randomness helps to decorrelate the trees and reduce overfitting

Q64 .How does a random forest reduce overfitting compared to a single decision tree?:

Random forests reduce overfitting compared to a single decision tree by averaging predictions from multiple trees trained on different subsets of data. By using bootstrap sampling and random feature selection, each tree in the random forest learns different aspects of the data, leading to less reliance on individual noisy or irrelevant features and improving generalization performance.

Q65 .What is bagging, and how is it related to random forests?:

Bagging (Bootstrap Aggregating) is a machine learning technique that involves training multiple models (often decision trees) independently on different subsets of the training data and then

combining their predictions through averaging or voting. Random forests are an extension of bagging specifically designed for decision trees, where each tree is trained on a bootstrap sample of the data and a random subset of features, resulting in a more diverse ensemble of trees.

Q66 .How do random forests handle categorical variables?:

Random forests handle categorical variables by splitting them into multiple binary (dummy) variables, where each category becomes a separate feature. During training, the algorithm considers these binary variables along with numerical variables for splitting nodes in the decision trees. This approach effectively incorporates categorical variables into the random forest model.

Q67.Can you explain the concept of feature importance in random forests?:

Feature importance in random forests measures the contribution of each feature to the predictive performance of the model. It is calculated based on the decrease in impurity (e.g., Gini impurity or entropy) achieved by each feature when used for splitting nodes in the decision trees. Features that result in larger impurity decreases are considered more important, as they provide more predictive power in distinguishing between classes or predicting the target variable.

Q68 .What are ensemble learning methods, and why are they used?:

Ensemble learning methods involve combining multiple individual models (learners) to improve predictive performance compared to any single model. They are used to mitigate the limitations of individual models, reduce variance, and enhance generalization by leveraging the wisdom of crowds. Ensemble methods often outperform single models by capturing diverse perspectives or patterns in the data.

Q69 .Explain the difference between bagging and boosting:

Bagging (Bootstrap Aggregating): Bagging involves training multiple models (often of the same type) independently on different subsets of the training data using bootstrap sampling. Predictions are then averaged (for regression) or majority-voted (for classification) to produce the final output.

Bagging reduces variance and overfitting by averaging the predictions of multiple models.

Boosting: Boosting is an iterative ensemble technique where models are trained sequentially, and each subsequent model focuses on correcting the errors made by the previous models. Boosting assigns higher weights to misclassified data points, thereby prioritizing difficult-to-predict instances. Examples of boosting algorithms include AdaBoost, Gradient Boosting, and XGBoost.

Q70 .What is the purpose of combining multiple weak learners in ensemble methods?:

The purpose of combining multiple weak learners (models that perform slightly better than random chance) in ensemble methods is to create a strong learner with improved predictive performance. By leveraging the diversity among weak learners and combining their predictions intelligently, ensemble methods can reduce bias, variance, and overfitting, leading to more robust and accurate models.

Q71 .How does stacking differ from bagging and boosting?:

Stacking (Stacked Generalization): Stacking involves training multiple heterogeneous (different types) or homogeneous (same type) base models and combining their predictions using a meta-learner (often a linear regression or neural network). Unlike bagging and boosting, which combine predictions in a parallel or sequential manner, stacking learns to weigh the predictions of base models optimally to minimize the overall error on a holdout dataset.

Q72 .Can you provide examples of ensemble methods other than random forests and boosting algorithms?:

Voting Classifiers: A voting classifier combines the predictions of multiple individual classifiers (e.g., decision trees, logistic regression, support vector machines) using a majority vote (for classification) or averaging (for regression).

Stacking: As mentioned earlier, stacking combines predictions from multiple base models using a meta-learner. It can involve various base models and meta-learner architectures, making it a versatile and powerful ensemble technique.

Q73.How do you evaluate the performance of decision trees and random forests?:

Performance of decision trees:

For classification: Metrics such as accuracy, precision, recall, F1-score, and ROC curve can be used.

For regression: Metrics like mean absolute error (MAE), mean squared error (MSE), and R-squared can be utilized.

Performance of random forests:

Similar to decision trees, but typically random forests are evaluated using metrics like out-of-bag (OOB) error rate, accuracy, and area under the ROC curve (AUC) for classification, and MSE or R-squared for regression.

Q74.What metrics can be used to measure the performance of ensemble models?:

Similar to individual models, ensemble models can be evaluated using metrics such as accuracy, precision, recall, F1-score, AUC, and mean squared error (MSE), depending on the type of problem (classification or regression) and the specific goals of the analysis.

Q75.Can you explain cross-validation and its role in evaluating ensemble models?:

Cross-validation involves splitting the dataset into multiple subsets (folds), training the model on some folds, and evaluating its performance on the remaining fold(s). This process is repeated multiple times, with different combinations of training and testing sets. Cross-validation helps to assess the model's generalization performance, detect overfitting, and estimate how the model will perform on unseen data. It is particularly useful for ensemble models as it provides more reliable estimates of performance by averaging results across multiple iterations.

Q77.What are some techniques for diagnosing and addressing overfitting in ensemble models?:

Techniques for addressing overfitting in ensemble models include:

1. Using simpler base models (e.g., shallow decision trees).
2. Regularization techniques such as limiting tree depth or node sample size.
3. Tuning hyperparameters (e.g., maximum tree depth, minimum samples per leaf) using cross-

- validation.
4. Feature selection or dimensionality reduction to reduce model complexity.
 5. Ensemble-specific techniques such as early stopping, where training is halted when performance on a validation set stops improving

Q79. Out-of-Bag (OOB) Error in Random Forests:

In a random forest, each decision tree is trained on a bootstrap sample of the original dataset, meaning that some data points are left out of each bootstrap sample. The out-of-bag (OOB) error is the prediction error calculated on the data points that are not included in the bootstrap sample used to train each individual tree. These out-of-bag data points serve as a validation set for the corresponding tree.

Q80 . decision tree Terminologies :

Decision Tree:

A decision tree is a supervised learning algorithm used for classification and regression tasks. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes.

Node:

A node in a decision tree represents a feature or attribute. There are two types of nodes:

Internal Node: Represents a decision point that splits the data into smaller subsets.

Leaf Node (or Terminal Node): Represents the final outcome or decision, usually a class label in the case of classification or a numerical value in the case of regression.

Root Node:

The root node is the topmost node in the decision tree, from which the tree starts to branch out. It represents the entire dataset.

Split:

A split is a decision point in a decision tree where the dataset is divided into two or more smaller subsets based on the value of a chosen attribute.

Branch:

A branch represents the outcome of a split and leads to subsequent nodes in the decision tree.

Pruning:

Pruning is a technique used to reduce the size of a decision tree by removing unnecessary branches.

This helps prevent overfitting and improves the tree's generalization ability.

Entropy:

As mentioned earlier, entropy is a measure of impurity or disorder in a set of data. It's used in decision tree algorithms to determine the best attribute to split the data on at each node.

Information Gain:

Information gain measures the effectiveness of an attribute in classifying the data. It quantifies the reduction in entropy or uncertainty after splitting the data on a particular attribute.

Gini Impurity:

Gini impurity is another measure of impurity used in decision tree algorithms. It measures the probability of misclassifying an element in a dataset, where a lower Gini impurity indicates purer nodes.

Decision Tree Pruning:

Decision tree pruning is a process of trimming down the branches of a decision tree to reduce its size and complexity. This helps to improve the tree's generalization performance on unseen data by reducing overfitting.

Q81 . What are some real-life applications of clustering algorithms?

The clustering technique can be used in multiple domains of data science like image classification, customer segmentation, and recommendation engine. One of the most common use is in market research and customer segmentation which is then utilized to target a particular market group to expand the businesses and profitable outcomes.

Q82 . What is feature engineering? How does it affect the model's performance?

Feature engineering refers to developing some new features by using existing features.

Sometimes there is a very subtle mathematical relation between some features which if

explored properly then the new features can be developed using those mathematical operations.

Also, there are times when multiple pieces of information are clubbed and provided as a single data column. At those times developing new features and using them help us to gain deeper insights into the data as well as if the features derived are significant enough helps to improve the model's performance a lot.

Q83. How do we measure the effectiveness of the clusters?

There are metrics like Inertia or Sum of Squared Errors (SSE), Silhouette Score, I₁, and I₂ scores. Out of all of these metrics, the Inertia or Sum of Squared Errors (SSE) and Silhouette score is a common metric for measuring the effectiveness of the clusters.

Q84.. Why do we take smaller values of the learning rate?

Smaller values of learning rate help the training process to converge more slowly and gradually toward the global optimum instead of fluctuating around it. This is because a smaller learning rate results in smaller updates to the model weights at each iteration, which can help to ensure that the updates are more precise and stable.

Q85. Why we cannot use linear regression for a classification task?

The main reason why we cannot use linear regression for a classification task is that the output of linear regression is continuous and unbounded, while classification requires discrete and bounded output values.

If we use linear regression for the classification task the error function graph will not be convex. A convex graph has only one minimum which is also known as the global minima but in the case of the non-convex graph, there are chances of our model getting stuck at some local minima which may not be the global minima. To avoid this situation of getting stuck at the local minima we do not use the linear regression algorithm for a classification task.

Q86. Why do we perform normalization?

To achieve stable and fast training of the model we use normalization techniques to bring all the features to a certain scale or range of values. If we do not perform normalization then there are chances that the gradient will not converge to the global or local minima and end up oscillating back and forth

Q87.What is the difference between upsampling and downsampling?

Upsampling: Increases the number of samples in the minority class by randomly selecting points and adding them to the dataset. It can lead to high training accuracy but might not generalize well to unseen data, as observed through lower validation accuracy.

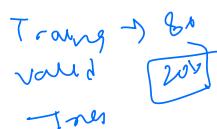
Downsampling: Decreases the number of samples in the majority class by randomly selecting a subset of points equal to the number of points in the minority class. It might result in the loss of some critical information due to data reduction.

Q88 .What is data leakage and how can we identify it?

Data leakage occurs when a high correlation exists between the target variable and input features.

Training the model with such highly correlated features can lead to the model inadvertently capturing the target variable's information during training. Consequently, the model may achieve seemingly high accuracy during training and validation, but it fails to perform well on new, unseen data.

Detecting this situation involves observing a significant drop in model performance when applied to real-world predictions.



Q89. Is it always necessary to use an 80:20 ratio for the train test split?

No there is no such necessary condition that the data must be split into 80:20 ratio. The main purpose of the splitting is to have some data which the model has not seen previously so, that we can evaluate the performance of the model.

Q90. What is one-shot learning?

One-shot learning is a concept in machine learning where the model is trained to recognize the

patterns in datasets from a single example instead of training on large datasets. This is useful when we haven't large datasets. It is applied to find the similarity and dissimilarities between the two images.

Q91 . What is the difference between one hot encoding and ordinal encoding?

One-Hot Encoding:

Represents each category with a binary vector.

Suitable for unordered categories without inherent order.

Creates binary features, increasing dimensionality.

Ideal for nominal variables like colors or genders.

Each category gets a separate binary feature.

Widely used in machine learning for categorical data.

Ordinal Encoding:

Assigns numerical values based on category order.

Suitable for ordered categories with a natural sequence.

Preserves order in a single feature, avoiding dimensionality increase.

Suitable for ordinal variables like education levels or ratings.

Numeric values are assigned based on the order of categories.

Helps capture the ordinal relationship between categories accurately.

Q92.How to identify whether the model has overfitted the training data or not?

This is the step where the splitting of the data into training and validation data proves to be a boon. If the model's performance on the training data is very high as compared to the performance on the validation data then we can say that the model has overfitted the training data by learning the patterns as well as the noise present in the dataset.

Q93.What is the difference between stochastic gradient descent (SGD) and gradient descent (GD)?

Gradient Descent (GD):

In GD, the model is trained on the entire dataset in each iteration.

It calculates the gradient of the loss function with respect to all training examples, averaging them to update the parameters.

The update step aims to move towards the global minimum of the loss function in a more consistent manner.

GD typically leads to smoother convergence of the training error towards the minimum.

Stochastic Gradient Descent (SGD):

In SGD, the model is trained using a randomly selected mini-batch of data in each iteration.

It calculates the gradient of the loss function with respect to only the examples in the mini-batch and updates the parameters based on this gradient.

The update steps are more noisy and less consistent compared to GD because they are influenced by a smaller subset of data.

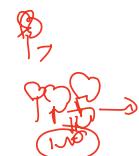
SGD often results in more oscillations in the training error, but it can converge faster and is computationally efficient, especially for large datasets.

$$K_{\text{min}} - kNN$$

Q94 . What is the difference between the k-means and k-means++ algorithms?

The only difference between the two is in the way centroids are initialized. In the k-means algorithm, the centroids are initialized randomly from the given points. There is a drawback in this method that sometimes this random initialization leads to non-optimized clusters due to maybe initialization of two clusters close to each other.

To overcome this problem k-means++ algorithm was formed. In k-means++, The first centroid is selected randomly from the data points. The selection of subsequent centroids is based on their separation from the initial centroids.



Q95. which one decision tree or random forest is more robust to outlier ?

Random Forest is generally more robust to outliers compared to a single Decision Tree.

Aggregation of Predictions: Random Forest aggregates predictions from multiple decision trees.

Outliers may have less influence on the final prediction since they are likely to be outliers in only a subset of the trees. This averaging effect helps to reduce the impact of outliers on the overall

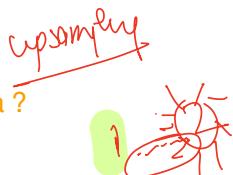
prediction.

Subsampling: Random Forest typically uses bootstrapping (sampling with replacement) to create multiple subsets of the data for each tree. This means that outliers may not always be present in every subset, further reducing their impact on individual trees' decisions.

Feature Randomization: Random Forest also randomly selects a subset of features at each split. This can help mitigate the effect of outliers in any single feature, as the model may rely on other features that are less affected by outliers.

Q96 How to handle imbalance data ?

Resampling Techniques:



Oversampling: Increase the number of instances in the minority class by generating synthetic samples or replicating existing ones.

Undersampling: Reduce the number of instances in the majority class by randomly selecting a subset of samples.

SMOTE (Synthetic Minority Over-sampling Technique): Generate synthetic samples for the minority class by interpolating between existing minority class instances.

Q97 .Explain SMOTE method used to handle data imbalance ?

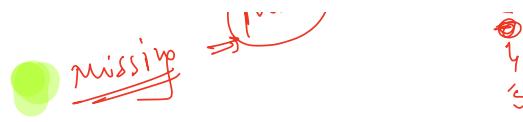
The Synthetic Minority Oversampling Technique (SMOTE) is employed to address data imbalance in datasets. This method synthesizes new data points for minority classes by linearly interpolating existing ones. Its advantage lies in training the model on augmented data, diversifying its exposure. However, SMOTE may introduce undesired noise, potentially harming model performance.

Q98.Does the accuracy score always a good metric to measure the performance of a classification model?

No, there are times when we train our model on an imbalanced dataset the accuracy score is not a good metric to measure the performance of the model. In such cases, we use precision and recall to measure the performance of a classification model.

(man)

2


Q99.What is KNN Imputer?

KNN Imputer is a sophisticated method for **filling null values in data**. It utilizes a distance parameter, also known as the **k parameter**, similar to clustering algorithms. Instead of using descriptive statistical measures like mean, mode, or median, KNN Imputer imputes missing values based on the neighborhood points of the missing values. This approach offers a more dynamic and context-sensitive imputation strategy.

Q100 .What is the purpose of splitting a given dataset into training and validation data?

Training Data: The training data is used to train the model's parameters or **learn patterns** from the data. The model learns the relationship between input features and the target variable by minimizing a loss function.

Validation Data: The validation data, also known as the holdout set, is used to **evaluate the model's performance during training**. It serves as an independent dataset to assess how well the model generalizes to unseen data. By evaluating the model on the validation set periodically during training, we can monitor its performance and make adjustments to hyperparameters or detect issues like overfitting.

Q101 . How can we visualize high-dimensional data in 2D?

One effective method is using the **t-SNE algorithm**, short for t-Distributed Stochastic Neighbor Embedding. It employs **non-linear complex techniques** to reduce data dimensionality. Alternatively, we can utilize **PCA or LDA** to convert n-dimensional data to 2D for visualization. However, while PCA aims to **preserve dataset variance**, t-SNE focuses on **retaining local similarities** in the dataset.

Q102 . Why do we face the curse of dimensionality?

The curse of dimensionality refers to the challenges and issues that arise when working with high-dimensional data. As the number of dimensions increases, the volume of the space grows exponentially, leading to various problems such as increased **computational complexity**, **sparsity of**

data, and the risk of overfitting. These issues make it difficult to effectively analyze, interpret, and model high-dimensional data, impacting the performance and reliability of machine learning algorithms.

Q103. What should be done if many columns are highly correlated?

Identify and Remove Redundant Features: Use techniques such as correlation analysis or variance inflation factor (VIF) to identify highly correlated features. Remove redundant features to reduce multicollinearity.

Feature Selection: Select a subset of features that are most relevant to the target variable. Consider techniques like recursive feature elimination, feature importance from tree-based models, or L1 regularization.

Dimensionality Reduction: Utilize dimensionality reduction techniques such as Principal Component Analysis (PCA) or t-distributed Stochastic Neighbor Embedding (t-SNE) to reduce the number of correlated features while preserving important information.

Feature Engineering: Create new features that capture the essence of the correlated features or transform existing features to reduce correlation. For example, you can create interaction terms or polynomial features.

Q104 . What is multicollinearity and how will you handle it in your regression model?

Multicollinearity refers to the presence of **high correlations between predictor variables** (also known as **independent variables**) in a regression model. It can cause issues such as unstable coefficient estimates and inflated standard errors, making the model's interpretation difficult.

→ input column

To handle multicollinearity in a regression model, several strategies can be employed:

Identify Multicollinearity: Use techniques such as **correlation analysis, variance inflation factor (VIF)**, or condition indices to identify highly correlated predictor variables.

Remove Redundant Variables: Remove one or more of the **highly correlated** variables from the model. Choose which variables to remove based on domain knowledge, importance to the model, or statistical significance.

Combine Variables: Create new composite variables by combining highly correlated variables. For example, you can calculate averages or sums of variables or create interaction terms.

Regularization Techniques: Use regularization techniques such as Ridge regression or Lasso regression. These methods penalize large coefficients, which can help mitigate multicollinearity issues by shrinking coefficients towards zero.



Q105 .What does the area under the ROC curve indicate?

ROC stands for **Receiver Operating Characteristic**. It measures the usefulness of a test where the larger the area, the more useful the test. These areas are used to compare the effectiveness of the tests. A higher **AUC (area under the curve)** generally indicates that the model is better at distinguishing between the positive and negative classes. AUC values range from 0 to 1, with a value of 0.5 indicating that the model is no better than random guessing, and a value of 1 indicating perfect classification.

Q106 what is sensitivity and specificity ?

Sensitivity (True Positive Rate):

$$\text{TP} = \frac{20}{100}$$

Sensitivity, also known as the true positive rate (**TPR**) or **recall**, measures the proportion of actual positive cases that are correctly identified by the model.

It is calculated as the ratio of true positives (TP) to the sum of true positives and false negatives (FN).

$$\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN})$$

Sensitivity indicates how well the model can identify positive instances from the total number of actual positive instances.

Specificity (True Negative Rate):

$$\frac{\text{TN}}{\text{TP} + \text{FN}}$$

Specificity measures the proportion of actual negative cases that are correctly identified by the model.

It is calculated as the ratio of true negatives (TN) to the sum of true negatives and false positives (FP).

$$\text{Specificity} = \text{TN} / (\text{TN} + \text{FP})$$

Specificity indicates how well the model can identify **negative instances** from the total number of

actual negative instances.

Q107 .What are some tools used to discover outliers?

Outliers are errors or extreme values that differ from the data in a set. They can impact the accuracy of your result if not detected through the outlier detection tools. Some popular tools to discover outliers include Z-score, Scatter plot, and Box plot.

Q108 . Explain Kernel SVM.

Kernel SVM stands for Kernel Support Vector Machine. In SVM, a kernel is a function that aids in problem-solving. They provide shortcuts to help you avoid doing complicated math. The amazing thing about kernel is that it allows us to go to higher dimensions and execute smooth computations. Kernel SVMs can work with a variety of kernel functions, including linear, polynomial, and radial basis function (RBF) kernels, among others.

Netflix, Amazon

Q109 . What is content-based filtering and collaborative filtering?

Content-based filtering recommends items to a user based on their past behavior, preferences, and interests. For instance, if a user has shown interest in action movies, the content-based filtering algorithm will recommend more action movies to that user.

Collaborative filtering recommends items to a user based on the preferences of other users who have similar tastes. For example, if a user has watched several action movies, the collaborative filtering algorithm will recommend other action movies that were also watched by users who have similar movie preferences.

Q110 .why does overfitting occur in machine learning ?

Overfitting happens when a model captures noise or random fluctuations in the data instead of the underlying pattern. It occurs due to:

Complex Models: Models with excessive complexity relative to the data can learn noise, causing overfitting.

Small Datasets: Insufficient data can lead the model to memorize examples rather than learning general patterns.

High Dimensionality: With many features, models may find false correlations that don't generalize.

Lack of Regularization: Inadequate regularization fails to penalize large parameter values, leading to overfitting.

Data Noise: Noise or outliers in the data can mislead the model into fitting irregularities.

Data Leakage: When test set information influences training, the model learns spurious patterns, causing overfitting.

Q111. 16. What is meant by Parametric and Non-parametric Models?

Parametric models refer to the models having a limited number of parameters. In case of parametric models, only the parameter of a model is needed to be known to make predictions regarding the new data.

Non-parametric models do not have any restrictions on the number of parameters, which makes new data predictions more flexible. In case of non-parametric models, the knowledge of model parameters and the state of the data needs to be known to make predictions.

Q112 .What is Entropy in Machine Learning?

Root \Rightarrow Gain

*Impurity < Entropy
Gini Index*

Entropy in Machine Learning measures the randomness in the data that needs to be processed. The more entropy in the given data, the more difficult it becomes to draw any useful conclusion from the data. For example, let us take the flipping of a coin. The result of this act is random as it does not favor heads or tails. Here, the result for any number of tosses cannot be predicted easily as there is no definite relationship between the action of flipping and the possible outcomes.

Q113 . When to use mean and when to use median to handle a missing numeric value?

We choose the mean to impute missing values when the data distribution is normal and there

are **no significant outliers**, as the mean is sensitive to both. In contrast, we use the median in cases of skewed distributions or when **outliers are present**, because the median is more robust to these factors and provides a better central tendency measure under these conditions.

Q114 . In Machine Learning, for how many classes can Logistic Regression be used?

Logistic regression cannot be used for **more than two classes**. Logistic regression is, by default, a binary classifier. However, in cases where **multi-class classification** problems need to be solved, the default number of classes can be extended, i.e., multinomial logistic regression.

Q115 What are the types of kernel in ML ?

Linear Kernel: Suitable for linearly separable data. It computes the dot product between the input features, often used for high-dimensional data.

Polynomial Kernel: Maps the data into higher-dimensional space using polynomial functions. It is effective for non-linear data and has parameters like degree and coefficient.

Radial Basis Function (RBF) Kernel: Measures the similarity between two data points based on their distance. It is versatile and widely used for various types of data, but it requires tuning of the gamma parameter.

Sigmoid Kernel: Applies a hyperbolic tangent function to the dot product of the input features. It is useful for neural networks and binary classification tasks.

Q116 Difference between StandardScaler and MinMaxScaler ?

StandardScaler (Z-score normalization):

Scales data to have a mean of 0 and a standard deviation of 1.

Suitable for algorithms assuming **normal distribution** and **robust to outliers**.

Transforms the data distribution to resemble a standard normal distribution.

MinMaxScaler:

Scales data to a specific range, often **between 0 and 1**.

Useful for models sensitive to feature magnitudes, but can be influenced by outliers.

$$\frac{x - \mu}{\sigma}$$

Preserves the original shape of the distribution, only changing the range of values.

Q117. How does one-hot encoding increase the dimensionality of a dataset compared to label encoding?

One-hot encoding increases the dimensionality of a dataset because it creates a separate binary variable (dummy variable) for each category in a categorical variable. For example, if a variable "Color" has three categories: "Red," "Green," and "Blue," one-hot encoding would create three binary variables: Color_Red, Color_Green, and Color_Blue.

label encoding assigns a unique integer value to each category without creating separate binary variables. For example, "Red" might be encoded as 0, "Green" as 1, and "Blue" as 2. Since label encoding does not create additional binary variables, it does not increase the dimensionality of the dataset.

☞ ↗ Numerical

Q118 . : What are the different types of Encoding techniques used in machine learning?

One-Hot Encoding: Creates binary variables for each category, with one variable set to 1 (hot) and the others set to 0 (cold). Each category is represented by a separate binary variable.

Label Encoding: Assigns a unique integer label to each category. Each category is mapped to a different integer value, allowing categorical variables to be represented numerically.

Ordinal Encoding: Similar to label encoding, but assigns integer labels based on the order or rank of the categories. Useful for ordinal variables where there is a natural ordering among categories.

Binary Encoding: Converts integer labels into binary digits and represents them as separate binary variables. It reduces dimensionality compared to one-hot encoding while capturing categorical information.

Frequency Encoding: Replaces categories with the frequency (count) of each category in the dataset. It captures the distribution of categories but may not work well for rare categories.

Target Encoding (Mean Encoding): Replaces categories with the mean of the target variable for each category. It captures the relationship between the categorical variable and the target

variable but may lead to overfitting if not used carefully.

Q119 .Various ways to do hyper parameter ?

Learning Rate: Controls the step size during optimization. It determines how much the model parameters are updated during each iteration of training.

Number of Trees (or Estimators): In ensemble methods like random forests and gradient boosting, this hyperparameter specifies the number of individual models (trees) to include in the ensemble.

Depth of Trees: Determines the maximum depth allowed for decision trees. Deeper trees can capture more complex relationships in the data but may also lead to overfitting.

Regularization Strength: Hyperparameters like alpha (L1 regularization) and lambda (L2 regularization) control the strength of regularization in models like linear regression and logistic regression. Regularization helps prevent overfitting by penalizing large parameter values.

Kernel Parameters: In kernel methods like Support Vector Machines (SVM) and kernelized versions of algorithms, hyperparameters like gamma control the shape and width of the kernel function.

Batch Size: Determines the number of training examples processed in each iteration of training in stochastic gradient descent and related optimization algorithms.

Activation Function: In neural networks, hyperparameters like the type of activation function (e.g., ReLU, sigmoid, tanh) and its parameters (e.g., slope in Leaky ReLU) affect the non-linear transformations applied to the data.

Number of Hidden Units/Layers: In neural networks, the number of hidden units (neurons) and layers determine the complexity and capacity of the model.

Dropout Rate: Controls the proportion of neurons randomly dropped out during training to prevent overfitting in neural networks.

Thresholds: In binary classification models, hyperparameters like the decision threshold affect

the classification of instances into positive or negative classes.

Q120. If your dataset is suffering from high variance, how would you handle it?

For datasets with high variance, we could use the **bagging algorithm** to handle it. Bagging algorithm splits the data into subgroups with sampling replicated from random data. After the data is split, random data is used to create rules using a training algorithm. Then we use polling technique to combine all the predicted outcomes of the model.

Q121 .What is a Box-Cox transformation?

Box-Cox transformation is a power transform which transforms **non-normal dependent variables into normal variables as normality** is the most common assumption made while using many statistical techniques. It has a **lambda parameter** which when set to 0 implies that this transform is equivalent to log-transform. It is used for **variance stabilization** and also to **normalize the distribution.**

Q122 .Advantages and Disadvantages of Decision Tree ?

Advantages:

Interpretability: Decision trees are **easy to understand and interpret**, making them suitable for explaining the logic behind the model's predictions to non-technical stakeholders.

Handle Non-linear Relationships: Decision trees can capture non-linear relationships between features and the target variable without the need for feature engineering or transformation.

Disadvantages of Decision Trees:

Overfitting: Decision trees are prone to overfitting, especially when the tree depth is not properly constrained or when dealing with noisy data. This can lead to poor generalization performance on unseen data.

Instability: Decision trees are sensitive to small variations in the training data, which can result in different trees being generated for slightly different datasets. This instability makes decision

trees less robust compared to some other machine learning algorithms.

Q123 .How do we check the normality of a data set or a feature?

Shapiro-Wilk W Test: This test assesses whether a sample comes from a normally distributed population. It is based on the correlation between the data and the corresponding normal distribution.

Anderson-Darling Test: The Anderson-Darling test evaluates whether a sample comes from a specified distribution, such as the normal distribution. It provides a single statistic that combines critical values from different percentiles.

Martinez-Iglewicz Test: The Martinez-Iglewicz test is used to detect outliers and assess normality based on the modified Z-score of the data points. It is particularly useful for detecting outliers in small to moderately sized datasets.



Q124 .What does the term Variance Inflation Factor mean?

Variation Inflation Factor (VIF) is the ratio of the model's variance to the model's variance with only one independent variable. VIF gives the estimate of the volume of multicollinearity in a set of many regression variables.

VIF = Variance of the model with one independent variable

Q125 .Which machine learning algorithm is known as the lazy learner, and why is it called so?

KNN is a Machine Learning algorithm known as a lazy learner. K-NN is a lazy learner because it doesn't learn any machine-learned values or variables from the training data but dynamically calculates distance every time it wants to classify, hence memorizing the training dataset instead.



Q126 .What is OOB error and how does it occur?

For each bootstrap sample, there is one-third of the data that was not used in the creation of the tree, i.e., it was out of the sample. This data is referred to as out of bag data. In order to

get an **unbiased measure** of the accuracy of the model over test data, **out of bag error** is used.

The out of bag data is passed for each tree is passed through that tree and the outputs are aggregated to give out of bag error.

Q127 .What is Bayes' Theorem? State at least 1 use case with respect to the machine learning context?

Bayes' Theorem describes the probability of an event, based on prior knowledge of conditions that might be related to the event. For example, if cancer is related to age, then, using Bayes' theorem, a person's age can be used to more accurately assess the probability that they have cancer than can be done without the knowledge of the person's age.

Chain rule for Bayesian probability can be used to predict the likelihood of the next word in the sentence.

Q128. Mention some of the EDA Techniques?

Exploratory Data Analysis (EDA) helps analysts to understand the data better and forms the foundation of better models.

Visualization

Univariate visualization

Bivariate visualization

Multivariate visualization

Missing Value Treatment – Replace missing values with Either Mean/Median

Outlier Detection – Use Boxplot to identify the distribution of Outliers, then **Apply IQR** to set the boundary for IQR

Transformation – Based on the distribution, apply a transformation on the features

Scaling the Dataset – Apply MinMax, Standard Scaler or Z Score Scaling mechanism to scale the data.

Feature Engineering – Need of the domain, and SME knowledge helps Analyst find derivative fields which can fetch more information about the nature of the data

Dimensionality reduction — Helps in reducing the volume of data without losing much information

Q129 . Difference between a generative and discriminative model

Generative Model: A generative model learns the **joint probability distribution** of the features and labels in the data. It models how the data is generated and can generate new samples similar to the training data. Generative models can be used for tasks such as data generation, missing data imputation, and semi-supervised learning.

Discriminative Model: A discriminative model learns the **conditional probability distribution** of the labels given the features. It focuses on learning the **boundaries or decision boundaries** between different classes in the data. Discriminative models directly model the decision boundary between classes and are typically used for tasks such as **classification** and **regression**.

Q130 . Which one is better, Naive Bayes Algorithm or Decision Trees?

Naive Bayes:

1. Work well with small dataset compared to DT which need more data
2. Lesser overfitting
3. Smaller in size and faster in processing

[Outlier, Less, Farter]

Decision Trees:

1. Decision Trees are very flexible, easy to understand, and easy to debug
2. No preprocessing or transformation of features required
3. Prone to overfitting but you can use pruning or Random forests to avoid that.

Q131 . What is log likelihood in logistic regression?

the log likelihood is a measure used to estimate how well the model fits the observed data. It

quantifies the agreement between the actual outcomes (labels) and the predictions made by the logistic regression model.

Q132 . Why does XGBoost perform better than SVM?

First reason is that XGBoosT is an ensemble method that uses many trees to make a decision so it gains power by repeating itself.

SVM is a linear separator, when data is not linearly separable SVM needs a Kernel to project the data into a space where it can separate it, there lies its greatest strength and weakness, by being able to project data into a high dimensional space SVM can find a linear separation for almost any data but at the same time it needs to use a Kernel and we can argue that there's not a perfect kernel for every dataset.

Q133 . What are the advantages of using a naive Bayes for classification?

1. Very simple, easy to implement and fast.
2. If the NB conditional independence assumption holds, then it will converge quicker than discriminative models like logistic regression.
3. Even if the NB assumption doesn't hold, it works great in practice.
4. Need less training data.
5. Highly scalable. It scales linearly with the number of predictors and data points.
6. Can be used for both binary and mult-i-class classification problems.
7. Can make probabilistic predictions.
8. Handles continuous and discrete data.
9. Not sensitive to irrelevant features.

