# Answers

**Question no. 1**

**Answer:**

**R-squared (R²):**

R-squared represents the proportion of the variance in the dependent variable that is predictable from the independent variables. It is a relative measure that explains how well the regression model fits the data. On the basis of Scale, it is easier to interpret. It allows for the comparison of different models.

**Residual Sum of Squares (RSS):**

RSS is the sum of the squared differences between the observed values and the values predicted by the model. It measures the total deviation of the response values from the fit to the response values. RSS is not standardized and its value depends on the scale of the dependent variable and the number of observations.

**In conclusion:**

R-squared is better measure of goodness of fit model in regression because it provides a clear and standardized measure of how well the model explains the variability in the dependent variable. It is easier to interpret and compare across different models and datasets. RSS, while useful, does not offer the same level of intuitive understanding and comparability.

**Question no. 2**

**Answer:**

In regression analysis, the Total Sum of Squares (TSS), Explained Sum of Squares (ESS), and Residual Sum of Squares (RSS) are key metrics used to evaluate the fit of a regression model.

**Total Sum of Squares (TSS):**

TSS measures the total variance in the observed data. It is the sum of the squared differences between the observed values and the mean of the observed values.

**Explained Sum of Squares (ESS):**

ESS measures the variance explained by the regression model. It is the sum of the squared differences between the predicted values and the mean of the observed values.

**Residual Sum of Squares (RSS):**

RSS measures the variance that is not explained by the regression model. It is the sum of the squared differences between the observed values and the predicted values.

**Relationship Between TSS, ESS, and RSS:**

The relationship between these three metrics is given by the equation:

**TSS = ESS+RSS**

This equation shows that the total variability in the data (TSS) is the sum of the variability explained by the model (ESS) and the variability that is not explained by the model (RSS).

**Question no. 3**

**Answer:**

**Regularization:** is a technique used in machine learning to prevent overfitting by adding a penalty to the model's complexity. It helps to address overfitting by adding a penalty for larger coefficients in the model. This encourages the model to keep the coefficients small and thus simpler.

Common regularization techniques are:

**1. L1 Regularization (Lasso):** The L1 regularization term adds the absolute value of the coefficients to the loss function.

**2. L2 Regularization (Ridge):** The L2 regularization term adds the squared value of the coefficients to the loss function.

**3. Elastic Net:**

Elastic Net combines both L1 and L2 penalties.

**Benefits of Regularization:**

**Improved Generalization**: Regularization helps models generalize better to new, unseen data by preventing overfitting.

**Simpler Models**: By penalizing large coefficients, regularization leads to simpler models that are easier to interpret.

**Feature Selection**: L1 regularization (Lasso) can perform automatic feature selection by shrinking some coefficients to zero.

**Complexity control:** we can control the complexity of the model and improve its performance on test data.

**Question no. 4**

**Answer:**

**The Gini impurity index**: is a crucial concept in decision tree algorithms, helping to measure the quality of splits and ensuring the construction of effective and accurate models.

It is a measure used in decision tree algorithms to evaluate the quality of a split. It quantifies the likelihood of a randomly chosen element being incorrectly classified if it were randomly labelled according to the distribution of labels in the subset.

Gini impurity ranges between 0 and 0.5 for binary classification and can be higher for multi-class classification.

- **A Gini impurity of 0:** indicates a perfectly pure node where all elements belong to a single class.

- **A Gini impurity close to 0.5:** (for binary classification) indicates a highly impure node where elements are evenly distributed across the classes.

In decision trees, the Gini impurity is used to decide the best split at each node. The algorithm calculates Gini impurity for all possible splits as well as chooses the split that results in the lowest weighted Gini impurity of the child nodes.

The goal is to achieve the largest reduction in impurity, thus ensuring that the splits created are as pure as possible, leading to more accurate and interpretable decision trees.

## Question no. 5

**Answer:**

Yes, unregularized decision trees are prone to overfitting. Here are the reasons why:

### 1. Complexity and Flexibility:

- **High Variance**: Decision trees are very flexible models. This flexibility can lead to high variance, where the model captures not only the underlying pattern but also the noise in the data.

- **Overly Complex Trees**: Without regularization, a decision tree can grow very deep, with many nodes and branches. Each additional split increases the likelihood of fitting noise in the training data rather than the true underlying pattern.

### 2. Perfectly Fit Training Data:

- **Pure Leaves**: An unregularized decision tree will keep splitting until all leaves are pure (i.e., all samples in a leaf belong to the same class) or until it perfectly fits the training data. This leads to a very specific model that performs well on the training data but poorly on new, unseen data.

- **High Depth and Many Leaves**: By growing very deep with many leaves, the model memorizes the training data instead of generalizing from it. This results in high training accuracy but poor test accuracy.

## 3. Lack of Pruning:

- **No Limitation on Growth**: Regularization techniques such as pruning limit the growth of the tree by setting constraints like maximum depth, minimum samples per leaf, or minimum samples required to split a node. Without these constraints, the tree grows excessively, capturing noise and leading to overfitting.

**Regularization Techniques to Prevent Overfitting**

Maximum Depth, Minimum Samples per Leaf, Minimum Samples per Split, Maximum Features, Pruning.

**Conclusion:**

Unregularized decision trees tend to overfit because they are highly flexible and can create complex models that capture noise in the training data. Applying regularization techniques helps control the complexity of the tree, ensuring better generalization and improved performance on unseen data.

## Question no. 6

**Answer:**

**Ensemble techniques:** in machine learning combine multiple models to produce a single, more powerful model. The primary goal of using ensemble methods is to improve the performance, robustness, and generalizability of the model compared to individual models. Key Aspects of Ensemble Techniques are Diversity and Aggregation.

Here are the types of ensemble techniques:

**Types of Ensemble Techniques:**

**1. Bagging (Bootstrap Aggregating)**: Multiple models are trained independently using different subsets of the training data, and their outputs are combined by averaging (for regression) or voting (for classification). For Example, Random Forests.

**2. Boosting**: Models are trained sequentially, with each new model focusing on the errors made by the previous models. The models are combined in a weighted manner. AdaBoost, Gradient Boosting Machines (GBM), XGBoost, LightGBM.

**3. Stacking (Stacked Generalization)**: Multiple base models (level-0 models) are trained, and their predictions are used as inputs to a higher-level model (meta-model), which makes the final prediction.

**4. Voting**: Combines predictions from multiple models by majority vote (for classification) or averaging (for regression).

**Hard Voting**: The class that gets the most votes is the final prediction.

**Soft Voting**: The probabilities of each class are averaged, and the class with the highest average probability is the final prediction.

**5. Bagging with Feature Selection (Random Subspaces)**: Similar to bagging, but also uses random subsets of features for training individual models, in addition to random subsets of data

**Benefits of Ensemble Techniques**

Improved Accuracy

Reduced Overfitting

Increased Robustness.

**In Conclusion,**

Ensemble techniques enhance machine learning model performance by leveraging the strengths and mitigating the weaknesses of individual models. By combining multiple models, these techniques improve accuracy, reduce overfitting, and increase robustness, making them powerful tools in the machine learning toolkit.

**Question no. 7**

**Answer:**

**Difference between Bagging and Boosting**

| Feature | Bagging | Boosting |
|---|---|---|
| **Training** | Parallel, independent training | Sequential, dependent training |
| **Data Sampling** | Bootstrap sampling (with replacement) | No bootstrap, weighted training |
| **Focus** | Reduces variance | Reduces bias |
| **Model Combination** | Averaging (regression), voting (classification) | Weighted combination |
| **Risk of Overfitting** | Generally lower | Higher, requires careful tuning |
| **Computational Efficiency** | Parallelizable | Less parallelizable, more sequential |
| **Examples** | Random Forests | AdaBoost: Gradient Boosting XGBoost, LightGBM |

**In Conclusion,**

Both bagging and boosting are powerful ensemble techniques but are suited to different types of problems and model characteristics. **Bagging** is effective at reducing variance and is well-suited for high-variance models, while **Boosting** is effective at reducing bias and

improving the performance of weak learners. The choice between them depends on the specific problem, the nature of the data, and the characteristics of the base learners being used.

## Question no. 8

**Answer:**

**Out-of-bag (OOB) error:** is an internal validation method used in random forests to estimate the prediction error without a separate validation set. Here's a brief overview:

1. **Bootstrap Sampling**: Each tree in a random forest is trained on a bootstrap sample (random sampling with replacement) of the original dataset.

2. **OOB Samples**: Data points not included in a tree's bootstrap sample are called OOB samples. On average, about 37% of the original data points are OOB for each tree.

3. **OOB Predictions**: Each tree makes predictions on its OOB samples. For each data point, predictions from all trees where it was OOB are aggregated.

4. **Error Estimation**: The OOB error is calculated by comparing these aggregated predictions to the actual values. For classification, this involves majority voting, and for regression, it involves averaging.

**Advantages:**

- **No Separate Validation Set Needed**: Utilizes all data for both training and validation.

- **Unbiased Estimate**: Provides an unbiased estimate of model performance.

The OOB error offers an efficient and reliable way to assess the accuracy of a random forest model.

**Question no. 9**

**Answer:**

**K-fold cross-validation** is a technique used to assess the performance and generalizability of a machine learning model. Here's a brief overview:

1. **Process**:

   - **Data Split**: The dataset is divided into $K$ equally sized folds (subsets).

   - **Training and Validation**: The model is trained $K$ times, each time using $K-1$ folds for training and the remaining fold for validation.

   - **Rotation**: This process is repeated such that each fold serves as the validation set exactly once.

2. **Error Estimation**:

   - **Aggregate Results**: The performance metric (e.g., accuracy, mean squared error) is calculated for each of the $K$ iterations.

   - **Average Performance**: The overall performance metric is obtained by averaging the results from all $K$ iterations.

**Advantages:**

- **Comprehensive Validation**: Uses all data points for both training and validation, providing a more reliable estimate of model performance.

- **Reduced Overfitting**: Helps in identifying how well the model generalizes to an independent dataset.

K-fold cross-validation is a robust method for model evaluation, ensuring that every data point is used for both training and validation, enhancing the reliability of the performance estimate.

**Question no. 10**

**Answer:**

Hyperparameter tuning in machine learning is the process of selecting the optimal values for hyperparameters, which are the parameters set before the learning process begins and control the model's training process and architecture.

**Why It Is Done**

1. **Optimize Performance:** Improve the model's accuracy and predictive power.

2. **Generalization:** Ensure the model performs well on unseen data by finding a balance between underfitting and overfitting.

3. **Control Complexity:** Manage the complexity of the model to prevent overfitting.

4. **Efficient Training:** Enhance training efficiency and speed by finding appropriate learning rates and batch sizes.

**Methods:**

1. Grid Search: Exhaustively searches over a predefined set of hyperparameters.

2. Random Search: Samples hyperparameters randomly from a distribution.

3. Bayesian Optimization: Uses probabilistic models to predict the performance of hyperparameter combinations.

4. Automated Machine Learning (Auto ML): Automated tools that perform hyperparameter tuning as part of an end-to-end pipeline.

Hyperparameter tuning is essential to maximize the performance and reliability of machine learning models.

# Question no. 11

**Answer:**

If the learning rate in Gradient Descent is too large, several issues can occur:

1. **Divergence**:

   o The algorithm may overshoot the minimum, causing the loss function to increase instead of decrease.

2. **Oscillations**:

   o The updates to the weights can be too drastic, leading to oscillations around the minimum rather than converging smoothly.

3. **Instability**:

   o The training process can become unstable, making it difficult for the model to converge to a solution.

4. **Suboptimal Convergence**:

   o Even if the algorithm converges, it might settle at a suboptimal point rather than the global minimum.

Using an appropriately small learning rate helps ensure stable and efficient convergence to the optimal solution.

# Question no. 12

**Answer:**

Logistic Regression is primarily suited for linear decision boundaries. For non-linear data, it may struggle to model the complexity of the relationships between features and classes.

**Limitations**:

- **Linear Decision Boundary**: Logistic Regression assumes a linear relationship between the features and the log-odds of the outcome.

- **Non-Linearity**: If the data is non-linear, Logistic Regression may not capture the complex patterns and interactions effectively.

**Workarounds**:

- **Feature Engineering**: Transforming or adding features to capture non-linearity.

- **Polynomial Features**: Adding polynomial terms to capture interactions.

- **Non-Linear Extensions**: Using kernel methods or more complex models like decision trees or neural networks.

In summary, while Logistic Regression is not inherently designed for non-linear classification, feature engineering or alternative models can be used to handle non-linear data.

## Question no. 13

**Answer:**

- **Adaboost**:
  - **Error Correction**: Sequentially adds models to correct the errors of previous models.
  - **Weight Adjustment**: Increases weights of misclassified data points to focus on difficult cases.
  - **Combining Models**: Uses weighted voting for combining model predictions.

- **Gradient Boosting**:

- o **Error Minimization**: Sequentially adds models to minimize residual errors using gradient descent.

- o **Flexible Loss Function**: Can optimize various loss functions (e.g., mean squared error, log loss).

- o **Model Combination**: Combines models by adding them to correct residual errors, not just misclassifications.

## Question no. 14

### Answer:

The bias-variance trade-off is the balance between two sources of error that affect the performance of a machine learning model:

- **Bias**: Error due to overly simplistic models that cannot capture the underlying patterns in the data. High bias leads to underfitting, where the model is too simple and fails to learn from the data effectively.

- **Variance**: Error due to overly complex models that capture noise in the training data as if it were a pattern. High variance leads to overfitting, where the model performs well on training data but poorly on unseen data.

**Trade-Off**:

- Increasing model complexity reduces bias but increases variance.

- Decreasing model complexity reduces variance but increases bias.

The goal is to find a balance where both bias and variance are minimized, leading to optimal model performance and generalization.

**Question no. 15**

**Answer:**

**Linear Kernel**: Computes the dot product of two vectors. Suitable for linearly separable data. The decision boundary is a straight line or hyperplane.

**RBF (Radial Basis Function) Kernel**: Measures similarity using a Gaussian function. Handles non-linear relationships by mapping data into a higher-dimensional space. The decision boundary can be highly flexible.

**Polynomial Kernel**: Computes a polynomial function of the dot product between vectors. Captures interactions between features. The decision boundary can represent polynomial relationships of various degrees.