# Assignment-5

## MACHINE LEARNING

**Question.1.** R-squared or Residual Sum of Squares (RSS) which one of these two is a better measure of goodness of fit model in regression and why?

**Answer.** R-squared is often preferred as a goodness-of-fit measure because it provides a standardized measure of the proportion of variance explained. It is easier to interpret and compare across different models.

RSS, on the other hand, directly measures the unexplained variance. It can be useful in situations where prediction accuracy is more critical than explaining variance.

In summary, R-squared is commonly used and provides a good overall assessment of model fit, but it's essential to consider both R-squared and RSS (or other metrics) together to get a comprehensive understanding of a regression model's performance.

**Question 2.** What are TSS (Total Sum of Squares), ESS (Explained Sum of Squares) and RSS (Residual Sum of Squares) in regression. Also mention the equation relating these three metrics with each other?

**Answer.** TSS represents the total variability in the dependent variable (Y) without considering any explanatory variables. It measures the squared differences between each observed Y value and the mean of Y.

ESS quantifies the variability in the dependent variable that is explained by the independent variables included in the regression model. It measures the squared differences between the predicted Y values and the mean of Y.

RSS represents the unexplained variability in the dependent variable, which is the part of the total variability not accounted for by the independent variables in the model. It measures the squared differences between the observed Y values and the predicted Y values.

The Equation relating these three metrics are:

TSS=ESS+RSS

**Question 3.** What is the need of regularization in machine learning?

**Answer.** Regularization is a technique used in machine learning to prevent overfitting and improve the generalization performance of a model. Overfitting occurs when a model learns the training data too well, capturing noise or random fluctuations that are not representative of the underlying patterns in the data. Regularization methods introduce additional

constraints or penalties on the model parameters during the training process to address this issue.

Common regularization techniques include L1 regularization (Lasso), L2 regularization (Ridge), and a combination of both (Elastic Net). The choice of regularization method and the strength of the regularization parameter are important hyperparameters that need to be tuned based on the specific characteristics of the data and the modeling task. Regularization is a valuable tool in the machine learning toolbox, particularly in situations where the risk of overfitting is high.

**Question 4.** What is Gini–impurity index?

**Answer.** The Gini impurity is a measure of impurity or uncertainty used in decision tree algorithms during the process of building a decision tree. It is commonly employed in classification problems to evaluate the quality of a split in the data based on the target variable.

The Gini impurity index for a particular node in a decision tree is calculated by assessing the probability of misclassifying a randomly chosen element in the dataset. In other words, it measures how often a randomly chosen element would be incorrectly classified based on the distribution of classes in the node.

**Question 5.** Are unregularized decision-trees prone to overfitting? If yes, why?

**Answer.** Yes, unregularized decision trees are prone to overfitting. Decision trees have the capacity to create highly complex structures that perfectly fit the training data, including capturing noise and outliers. Overfitting occurs when a model learns the training data too well, to the point that it captures the noise and specific patterns that are present only in the training set but do not generalize well to new, unseen data.

**Question 6.** What is an ensemble technique in machine learning?

**Answer**.  An ensemble technique in machine learning involves combining multiple individual models to create a stronger and more robust model. The idea behind ensemble methods is to leverage the diversity of different models and combine their predictions to improve overall performance, increase accuracy, and reduce the risk of overfitting.

**Question 7.** What is the difference between Bagging and Boosting techniques?

**Answer. Training Approach:**

    1.Bagging trains models independently in parallel.

    2.Boostin trains models sequentially,with each model focused on correcting the errors

of the previous models.

Data Sampling:

1.Bagging involves creating multiple subsets of the data through bootstrap sampling.

2.Boosting adjusts weights of instance during training to focus on misclassified.

In summary,While both Bagging and boosting are ensemble technique that aim to improve model performance, they differ in their approach to data sampling, training, and combining predictions. Bagging aims to reduce variance and increase stability, while Boosting focuses on sequentially building models that correct errors, with an emphasis on improving accuracy.

**Question 8.**What is out-of-bag error in random forests?

**Answer.** The out-of-bag error provides an unbiased estimate of the model's performance on new, unseen data, as the data points used for estimation were not part of the training set for the particular tree being evaluated. This estimation is useful for assessing the generalization ability of the Random Forest model and can be used as an alternative or complement to cross-validation.

**Quesiton 9.** What is K-fold cross-validation?

**Answer.** K-fold cross-validation is a resampling technique used in machine learning to assess the performance and generalization ability of a model. The dataset is divided into K subsets (or folds), and the model is trained and evaluated K times, each time using a different fold as the test set and the remaining folds as the training set. This process allows for a more robust estimation of the model's performance compared to a single train-test split.

K-fold cross-validation is particularly useful for hyperparameter tuning, model selection, and comparing the performance of different algorithms. It provides a more robust and unbiased estimate of a model's ability to generalize to new, unseen data compared to a single train-test split.

**Question 10.** What is hyper parameter tuning in machine learning and why it is done?

**Answer.** Hyperparameter tuning, also known as model optimization or hyperparameter optimization, is the process of finding the best set of hyperparameters for a machine learning model. Hyperparameters are configuration settings external to the model that cannot be learned from the training data but significantly impact the model's performance.

Examples of hyperparameters include the learning rate in gradient boosting, the number of hidden layers and neurons in a neural network, the regularization term in linear regression, and the depth of decision trees in a random forest.

There are various techniques for hyperparameter tuning, including grid search, random search, and more advanced methods like Bayesian optimization. The process involves defining a search space for hyperparameters, selecting a search strategy, training and evaluating models with different hyperparameter configurations, and selecting the combination that yields the best performance.

In summary, hyperparameter tuning is a crucial step in the machine learning model development process, helping to optimize model performance, prevent overfitting or underfitting, and improve the generalization of the model to new, unseen data.

**Question 11.** What issues can occur if we have a large learning rate in Gradient Descent?

**Answer.** A large learning rate in gradient descent can lead to several issues, and it can adversely affect the training of machine learning models. Here are some common issues associated with a large learning rate:

**a .Divergence:**

One of the most significant problems with a large learning rate is that it can cause the optimization process to diverge. Instead of converging to the minimum of the loss function, the algorithm may oscillate or diverge, leading to increasingly large updates in parameter values.

**b. Overshooting the minimum:**

A large learning rate can cause the algorithm to overshoot the minimum of the loss function. This means that the updates to the model parameters are too large, and the algorithm may oscillate around the minimum without converging.

**c. Instability and unstable gradient:**

Large learning rates can lead to numerical instability, especially in the presence of high curvature or steep gradients in the loss function. This instability can make it challenging to find a stable and accurate solution.

**d. Missed Convergence:**

With a large learning rate, the algorithm might skip over the optimal solution, missing the convergence point. This can result in a suboptimal or even an entirely incorrect model.

**e. .Slow Convergence in the long run:**

Paradoxically, while a large learning rate may cause rapid updates initially, it can slow down the convergence in the long run. The algorithm may oscillate back and forth, taking longer to reach convergence or failing to converge at all.

**f. Increased risk for overfitting:**

Large learning rates can lead to overfitting by causing the model to memorize the training data, capturing noise rather than learning the underlying patterns. This is especially problematic when training deep neural networks.

**Question 12.** Can we use Logistic Regression for classification of Non-Linear Data? If not, why?

**Answer.** Logistic Regression is a linear model designed for binary classification tasks. While it's a powerful and widely used algorithm, its linearity assumption means that it works well when the decision boundary between classes is approximately linear. Logistic Regression models the relationship between the input features and the log-odds of belonging to a particular class, and it uses a logistic function (sigmoid) to transform this into probabilities.

However, Logistic Regression may not perform well on non-linear data because it is limited to capturing linear relationships. If the true decision boundary in the data is non-linear, Logistic Regression may struggle to accurately model and classify the data.

In cases where the decision boundary is non-linear, using a linear model like Logistic Regression may lead to underfitting, where the model is too simplistic to capture the underlying patterns in the data. This limitation is evident in situations where the classes are not separable by a straight line in the feature space.

**Question 13.** Differentiate between Adaboost and Gradient Boosting?

**Answer**. Adaboost (Adaptive Boosting) and Gradient Boosting are both ensemble learning techniques, but they differ in their approach to building a strong predictive model through the combination of weak learners (typically decision trees). Here are the key differences between Adaboost and Gradient Boosting:

    a. **Loss function:**

**Adaboost:** Adaboost focuses on minimizing the classification error. It assigns higher weights to misclassified data points, and subsequent weak learners are trained to emphasize correcting the errors made by the previous ones.

**Gradient Boosting:** Gradient Boosting minimizes a specified loss function, which can be customized based on the problem (e.g., mean squared error for regression, log loss for classification). It sequentially fits models to the residuals (or negative gradients) of the loss function.

    b. **Weaking of weak learners:**

**Adaboost:** It assigns different weights to each weak learner based on its performance. A higher weight is given to the weak learners that perform well on the training data.

**Gradient Boosting:** Weak learners are added sequentially, with each one correcting the errors of the combined ensemble. The contribution of each weak learner is determined by the gradients of the loss function.

### c. Model Training:

**Adaboost:** The algorithm sequentially fits a series of weak learners to the data, adjusting the weights of misclassified instances at each step. The final model is a weighted sum of the weak learners.

**Gradient Boosting:** The algorithm sequentially builds a series of weak learners (usually decision trees) by fitting them to the residuals of the previous models. Each new tree corrects the errors made by the combined ensemble of previous trees.

### d. Combination prediction:

**Adaboost:** The final prediction is a weighted sum of the weak learners, where weights are determined by their performance.

**Gradient Boosting:** The final prediction is the sum of the predictions made by all the weak learners, each multiplied by a learning rate.

### e. Learning rate:

**Adaboost:** Typically uses a learning rate of 1.0. Adjustments are made through the weighting of weak learners.

**Gradient Boosting:** Uses a learning rate (shrinkage parameter) to control the contribution of each weak learner. A lower learning rate requires more weak learners but may lead to better generalization.

### f. Parallelism:

**Adaboost:** Can be parallelized as the weak learners are independent of each other.

**Gradient Boosting:** Generally harder to parallelize due to the sequential nature of training weak learners.

**Question 14.** What is bias-variance trade off in machine learning?

**Answer.** **1. Bias**

**Definition:** Bias refers to the error introduced by approximating a real-world problem, which may be complex, by a simplified model. It represents the difference between the predicted values and the true values.

**Characteristics:** High bias often leads to underfitting, where the model is too simplistic and fails to capture the underlying patterns in the data. The model is not able to learn from the training data effectively.

2. Variance:

**Definition:** Variance is the amount by which the model's predictions would change if it were trained on a different dataset. It measures the model's sensitivity to the variability in the training data.

**Characteristics:** High variance can lead to overfitting, where the model fits the training data too closely, capturing noise and not generalizing well to new, unseen data.

3. Trade off:

**Bias-Variance Trade-off:** The bias-variance trade-off implies that there is a trade-off between model bias and variance. A model with high bias tends to have low variance, and vice versa. The goal is to find the right balance that minimizes the total error on new, unseen data.

**4. Total Error:**

Irreducible Error: This is the error that cannot be reduced, as it is inherent in the problem itself. It may be due to noise in the data or other factors that cannot be controlled.

The trade-off aims to minimize the sum of bias$^2$ and variance to achieve the lowest possible total error.

**Question 15.** Give short description each of Linear, RBF, Polynomial kernels used in SVM.
**Answer.**

1. **Linear Kernel:**
   **Definition:**The linear kernel is the simplest kernel and is used for linearly separable data. It represents the dot product of the input features and is effective when the relationship between the features and the target variable is approximately linear.
   **Use Case:** Suitable for linearly separable data or cases where a linear decision boundary is appropriate.

2. **Radial Basis Function:**
   **Definition:** The RBF kernel, also known as the Gaussian kernel, transforms data into an infinite-dimensional space. It is characterized by a parameter (gamma) that determines the influence of each training example on the decision boundary. The RBF kernel is versatile and can capture complex non-linear relationships in the data.

   **Use Case:** Suitable for non-linear data and situations where a flexible decision boundary is needed. However, it requires careful tuning of the $(Gamma)$ parameter.

3. **Polynomial Kernel:**

   **Definition:** The polynomial kernel allows SVM to handle polynomial relationships by transforming the input features into higher-dimensional spaces. It is defined by a degree parameter $d$, which determines the degree of the polynomial. Higher degrees increase the complexity of the decision boundary.

   **Use Case:** Suitable for capturing polynomial relationships in the data. The degree parameter should be chosen based on the complexity of the underlying patterns in the data.