

Name: Adeola Odunewu  
Internship N0: DS2311  
Assessment 5: PFA of Worksheet Set 2

## Machine Learning Worksheet Solutions

1.

The choice between R-squared ( $R^2$ ) and Residual Sum of Squares (RSS) depends on the specific goals and priorities of your analysis. There is no clear answer as to which one is universally "best" because each metric provides different insights into model performance:

### R-squared ( $R^2$ ):

**Strengths:** Provides a clear interpretation of the proportion of variance in the dependent variable explained by the model. A higher R-squared suggests a better fit.

**Considerations:** Can increase even with the addition of irrelevant variables, and it doesn't penalize for overfitting.

### Residual Sum of Squares (RSS):

**Strengths:** Directly measures the overall magnitude of prediction errors, offering insight into how well the model predicts the observed data.

**Considerations:** Does not provide a standardized measure, and it does not account for the scale of the dependent variable.

2.

In regression analysis, Total Sum of Squares (TSS), Explained Sum of Squares (ESS), and Residual Sum of Squares (RSS) are terms used to quantify the variability in the dependent variable (response variable).

TSS represents the total variability in the dependent variable (Y) without considering the effects of the independent variables (X). It is essentially the sum of the squared differences between each observed dependent variable value and the mean of the dependent variable.

$TSS = \sum (Y_i - \bar{Y})^2$   $Y_i$  is each observed value of the dependent variable, and  $\bar{Y}$  is the mean of the dependent variable.

ESS measures the variability in the dependent variable that is explained by the independent variables (X) in the regression model. It represents the sum of the squared differences between the predicted values (based on the regression model) and the mean of the dependent variable.

$ESS = \sum (\hat{Y}_i - Y_i)^2$  where  $\hat{Y}$  is the predicted value for each observation.

RSS quantifies the unexplained variability in the dependent variable. It is the sum of the squared differences between the observed values of the dependent variable and the predicted values from the regression model.

$ESS = \sum (Y_i - \hat{Y}_i)^2$  Where  $Y_i$  is each observed value, and  $\hat{Y}_i$  is the predicted value.

The relationship between these three metrics is given as:

$$TSS = ESS + RSS$$

This equation illustrates that the total variability in the dependent variable (TSS) can be decomposed into two components: the variability explained by the regression model (ESS) and the unexplained residual variability (RSS). In a well-fitted model, ESS should be relatively high compared to RSS, resulting in a higher  $R^2$  (coefficient of determination), where  $R^2 = TSS/ESS$

### 3.

Regularization in machine learning is essential to address challenges such as overfitting, especially in high-dimensional datasets with many features. It adds penalty terms to the model's objective function, discouraging overly complex models and preventing them from fitting noise in the training data. Regularization techniques, including L1 (Lasso) and L2 (Ridge) regularization, are particularly useful in handling collinearity among features and improving the generalization performance of models. Additionally, regularization aids in ensuring numerical stability during the optimization process, promoting more reliable and robust model training. Overall, regularization is a crucial tool to strike a balance between fitting the training data well and preventing models from becoming too complex and overfitting.

### 4.

The Gini impurity index is a measure commonly used in decision tree algorithms to evaluate the impurity or disorder of a set of data points within a node. In the context of decision trees, a node represents a subset of the data, and the Gini impurity is used to assess how often a randomly chosen element would be incorrectly classified if it were randomly assigned a class label according to the distribution of classes in the node.

$$Gini(t) = 1 - \sum_{I=1}^C P(I|t)$$

Where  $C$  is number of class

$P(I|t)$  is the probability of class  $I$  in node  $t$ .

### 5.

**Yes.** Unregularized decision trees are prone to overfitting, which occurs when the model fits the training data too closely and captures noise or specific patterns that do not generalize well. Decision trees, being flexible, can become overly complex,

memorize noise, and be sensitive to small variations in the training data. To address overfitting, regularization techniques are applied, including pruning (removing unnecessary branches), limiting tree depth, and imposing constraints on node size. These regularization methods help strike a balance between capturing important patterns and preventing the model from fitting the training data too closely, resulting in better generalization to new, unseen data.

## 6.

Ensemble techniques in machine learning involve combining predictions from multiple individual models to enhance overall performance. Two common types are:

**Bagging (Bootstrap Aggregating):** \*\* Trains multiple instances of the same algorithm on different subsets of the data, and their predictions are averaged or majority-voted. Example: Random Forest, which employs multiple decision trees.

**Boosting:** Sequentially trains weak learners, giving more weight to misclassified instances from previous models, leading to improved correction of errors. Example Algorithms: AdaBoost, Gradient Boosting, and XGBoost.

Ensemble methods offer advantages such as improved generalization, increased accuracy, robustness to noise, and versatility across different machine learning algorithms. They are widely used and have contributed to the success of various applications.

## 7.

Bagging (Bootstrap Aggregating) and Boosting are both ensemble techniques in machine learning, but they differ in their approaches. Bagging trains multiple models independently on different subsets of the data, with predictions combined through averaging or voting. It is robust to outliers and noise.

Boosting sequentially trains models, giving more weight to instances misclassified by previous models. Predictions are combined with weighted averaging. Boosting is sensitive to outliers and aims to reduce bias.

Both techniques improve model performance by leveraging multiple models, but they have distinct training processes, handling of instances, and ways of combining predictions. Common algorithms include Random Forest for Bagging and AdaBoost, Gradient Boosting, and XGBoost for Boosting.

## 8.

The out-of-bag (OOB) error in Random Forests is a measure of the model's performance on data points that were not included in the bootstrap samples used for

training each tree. For each tree, the OOB data points (not included in its bootstrap sample) are used to evaluate the model's prediction performance.

The OOB errors from all trees are then averaged or combined to obtain an overall estimate of the Random Forest's generalization error. The OOB error serves as an unbiased indicator of how well the model is likely to perform on new, unseen data and is useful for model tuning and assessment without the need for a separate validation set.

#### **9.**

K-fold cross-validation is a technique for evaluating machine learning model performance. It involves splitting the dataset into K equally sized folds, training and validating the model K times, with each fold used as the validation set exactly once. This process helps reduce variance, ensures better use of data, and provides a robust estimate of the model's generalization performance by averaging performance metrics across the K iterations. Common choices for K include 5-fold and 10-fold cross-validation, making it a valuable tool for assessing model performance and tuning hyperparameters.

#### **10.**

Hyperparameter tuning in machine learning involves optimizing external configuration settings (hyperparameters) of a model to improve its performance and generalization. These settings, not learned from data, impact the model's behavior. The goals of hyperparameter tuning include enhancing model performance, preventing overfitting, achieving efficient training, and ensuring good generalization to new data.

#### **11.**

Using a large learning rate in gradient descent can lead to various issues like:

Divergence, Overshooting the Minimum, Instability, Failure to Converge, Inability to Generalize

#### **12.**

Logistic Regression is inherently a linear model, and its decision boundary is a linear function. Therefore, when faced with non-linear data, Logistic Regression may not perform well as it cannot capture complex non-linear relationships between features and the target variable. If the decision boundary in the data is non-linear, Logistic Regression might struggle to accurately classify instances.

#### **13.**

While both Adaboost and Gradient Boosting are ensemble techniques, they differ in their approach to assigning weights to data points, the loss function they minimize, and how subsequent models contribute to the final prediction.

Adaboost focuses on adjusting weights to emphasize misclassified points.

Gradient Boosting minimizes the residuals through sequential optimization.

#### 14.

The bias-variance tradeoff is a key concept in machine learning, describing the balance between bias and variance in model performance. Bias represents the error from oversimplified models, leading to underfitting, while variance represents sensitivity to noise, leading to overfitting. There is a tradeoff, as decreasing bias often increases variance and vice versa. Striking the right balance is crucial for optimal model generalization. High bias models are too simple, high variance models are too complex, and the goal is to find an optimal level of complexity.

#### 15.

##### i. Linear Kernel:

**Description:** The linear kernel is the simplest form of the SVM kernel and is used for linearly separable data. It computes the dot product between the feature vectors in the original input space.

**Use Case:** Suitable for linearly separable datasets where the decision boundary is a hyperplane.

##### ii. Radial Basis Function (RBF) Kernel:

**Description:** The RBF kernel, also known as the Gaussian kernel, transforms the input features into an infinite-dimensional space. It measures the similarity between two examples based on the Euclidean distance.

**Use Case:** Effective for capturing complex, non-linear relationships in data. It is widely used when the decision boundary is expected to be non-linear.

##### iii. Polynomial Kernel:

**Description:** The polynomial kernel computes the dot product of the feature vectors raised to a specified power, introducing non-linearity to the model. The degree parameter controls the degree of the polynomial.

**Use Case:** Useful for capturing polynomial relationships between features. It can model curved decision boundaries.