

## MACHINE LEARNING WORKSHEET ANSWERS :

### Q1]

RSS is more sensitive to outliers:  $R^2$  can be influenced by a single outlier, which can artificially inflate the value. RSS, on the other hand, is more robust to outliers, as it penalizes the model for large residuals.

RSS provides a more absolute measure:  $R^2$  is a relative measure that depends on the scale of the data. RSS, being an absolute measure, provides a more direct assessment of the model's performance.

RSS is more closely related to the model's predictive power: A lower RSS indicates that the model is able to make more accurate predictions, which is ultimately the goal of regression analysis.

In summary, while  $R^2$  provides a useful measure of the model's explanatory power, RSS is a more informative metric for evaluating the goodness of fit of a regression model.

### Q2]

TSS (Total Sum of Squares): TSS measures the total variation in the dependent variable ( $y$ ). It represents the sum of the squared differences between each observed value and the mean of the dependent variable.

ESS (Explained Sum of Squares): ESS measures the variation in the dependent variable that is explained by the independent variables in the model. It represents the sum of the squared differences between each predicted value and the mean of the dependent variable.

RSS (Residual Sum of Squares): RSS measures the variation in the dependent variable that is not explained by the independent variables in the model. It represents the sum of the squared differences between each observed value and the predicted value.

The equation relating these three metrics is:

$$TSS = ESS + RSS$$

This equation indicates that the total variation in the dependent variable (TSS) can be decomposed into two components: the variation explained by the model (ESS) and the variation not explained by the model (RSS). Additionally, the coefficient of determination ( $R^2$ ) can be calculated using these metrics:

$$R^2 = ESS / TSS \text{ Or } R^2 = 1 - (RSS / TSS)$$

$R^2$  measures the proportion of the total variation in the dependent variable that is explained by the independent variables in the model.

### **Q3]**

The need for regularization arises from the fact that complex models with many parameters can easily fit the noise in the training data, rather than the underlying patterns. This can lead to overfitting, which can be mitigated by adding a penalty term to the loss function that discourages large parameter values.

Regularization techniques, such as L1 and L2 regularization, add a penalty term to the loss function that encourages the model to have smaller parameter values. This helps to reduce the complexity of the model and prevent overfitting.

L1 regularization adds a penalty term proportional to the absolute value of the parameters, which can lead to sparse solutions where some parameters are set to zero. L2 regularization adds a penalty term proportional to the square of the parameters, which tends to shrink the parameter values towards zero but does not set them to zero.

Regularization can also help to improve the interpretability of the model, as smaller parameter values can make the model more interpretable and easier to understand.

In summary, the need for regularization in machine learning arises from the risk of overfitting, which can lead to poor generalization to new data. Regularization techniques can help to prevent overfitting by adding a penalty term to the loss function that encourages the model to have smaller parameter values, leading to improved performance on new data and increased interpretability.

### **Q4]**

The Gini-impurity index is a measure of the impurity or disorder in a dataset. It is used in decision trees to determine the best attribute to split on at each node. The Gini-impurity index is calculated as 1 minus the sum of the squared probabilities of each class in the dataset. It varies between 0 and 1, with 0 indicating a pure dataset (all elements belong to a single class) and 1 indicating a completely impure dataset (elements are randomly distributed across various classes). The goal of decision tree algorithms is to minimize the Gini-impurity index, and thus the amount of impurity, at each node.

The formula for calculating the Gini-impurity index is:

$$Gini_{\{A\}}(D) = \frac{n_1}{n} Gini(D_1) + \frac{n_2}{n} Gini(D_2)$$

where  $n$  is the total number of samples in the dataset  $D$ ,  $n_1$  and  $n_2$  are the number of samples in the subsets  $D_1$  and  $D_2$ , respectively, and  $Gini(D_1)$  and  $Gini(D_2)$  are the Gini-impurity indices of the subsets  $D_1$  and  $D_2$ , respectively.

## Q5]

Yes, unregularized decision trees are prone to overfitting. This is because they can learn the training set to a high degree of granularity, making them easily overfit. Allowing a decision tree to split to a granular degree enables it to learn every point extremely well, resulting in perfect classification on the training set, but poor generalization on new, unseen data.

In other words, decision trees can become too specialized in fitting the training data, which leads to poor performance on test data. This is especially true when the decision tree is allowed to grow without any constraints, resulting in a complex model that is highly tailored to the training data.

To avoid overfitting, it's essential to implement regularization techniques, such as pruning, setting a maximum depth, or using ensemble methods like bagging or random forests. Additionally, ensuring that the training and test sets are representative of the problem domain and have enough data to generalize well can also help mitigate overfitting.

## Q6]

An ensemble technique in machine learning is a method that combines the predictions from multiple models to achieve better predictive performance. This can be done through various approaches, such as bagging, boosting, and stacking.

Bagging involves creating multiple subsets of the original dataset and training a model on each subset. The predictions from each model are then combined using a voting or averaging scheme to produce the final prediction. Boosting involves training a sequence of models, with each subsequent model being trained on the residual errors of the previous model.

The final prediction is then made by combining the predictions from each model. Stacking involves training multiple models on the same dataset and then combining their predictions using another machine learning model, known as a meta-learner or second-level model.

Ensemble techniques can help to reduce overfitting, improve accuracy, and provide more robust predictions compared to using a single model.

## **Q7]**

Bagging and Boosting are both ensemble learning techniques used in machine learning to improve the accuracy and stability of machine learning algorithms. However, they differ in how they combine the predictions from multiple models.

Bagging, or Bootstrap Aggregating, involves creating multiple subsets of the original dataset and training a model on each subset. The predictions from each model are then combined using a voting or averaging scheme to produce the final prediction. Bagging helps to reduce the variance and avoid overfitting of the model. It is usually applied to decision tree methods.

Boosting, on the other hand, involves training a sequence of models, with each subsequent model being trained on the residual errors of the previous model. The final prediction is then made by combining the predictions from each model. Boosting helps to reduce the bias and improve the accuracy of the model. It is usually applied to weak models, such as decision stumps or shallow decision trees.

Bagging is a parallel ensemble technique that reduces variance and helps to avoid overfitting, while Boosting is a sequential ensemble technique that reduces bias and improves accuracy.

## **Q8]**

Out-of-bag (OOB) error is a method used to measure the prediction error of random forests, a type of ensemble learning algorithm. It is calculated by leaving out a subset of the training data when building each tree in the forest, and then using that left-out data to test the predictions of the tree. The OOB error is then calculated as the average error across all the trees in the forest. This approach provides a reliable estimate of the generalization error of the random forest model, without the need for a separate validation dataset. The OOB error can be used for model selection, tuning, and evaluation, and it is particularly useful for identifying the optimal number of trees in the forest.

### Q9]

K-fold cross-validation is a resampling procedure used to evaluate machine learning models on a limited data sample. The basic idea is to split the available data into  $k$  smaller subsets, called folds. Then, for each fold, the model is trained on the remaining  $k-1$  folds and evaluated on the current fold. This process is repeated  $k$  times, and the average performance across all folds is used as an estimate of the model's performance.

### Q10]

Hyperparameter tuning is the process of finding the optimal hyperparameters for a machine learning model that results in the best performance on a given dataset. Hyperparameters are parameters that are set before training a model, such as learning rate, batch size, number of hidden layers, and regularization strength. These hyperparameters cannot be learned from the data and must be specified by the user.

Hyperparameter tuning is done to improve the performance of a machine learning model. Different hyperparameters can significantly affect the model's performance, and finding the optimal combination of hyperparameters can be a challenging task. Hyperparameter tuning helps to:

- Improve model accuracy
- Reduce overfitting or underfitting
- Increase model interpretability
- Enhance model robustness

### Q11]

If the learning rate in Gradient Descent is too large, the model may converge too quickly to a suboptimal solution. This is because a large learning rate may cause the model to overshoot the optimal point during the gradient descent process, resulting in unstable training. Additionally, a learning rate that is too large may cause the optimization process to diverge or oscillate, making it difficult to find the optimal weights

for the model. Therefore, it is important to carefully select the learning rate to ensure that the model converges to the optimal solution

## **Q12]**

No, we cannot use Logistic Regression for classification of Non-Linear Data. Logistic Regression constructs linear boundaries and assumes a linear relationship between the dependent variable and the independent variables. It is not suitable for non-linear problems because it has a linear decision surface. In real-world scenarios, linearly separable data is rarely found, and Logistic Regression is not capable of handling non-linear relationships.

However, there are some techniques to extend Logistic Regression to handle non-linear relationships. One approach is to use polynomial features, where you create new features by taking the polynomial combinations of the original features. Another approach is to use kernel methods, which can be used to transform the data into a higher-dimensional space where the relationship becomes linear.

Additionally, you can also use other algorithms such as Decision Trees, Random Forest, Support Vector Machines (SVM), or Neural Networks, which are more suitable for handling non-linear relationships.

It's also worth noting that Logistic Regression can be used for non-linear classification if the non-linearity is introduced through the use of non-linear features or transformations. For example, you can use a polynomial or radial basis function (RBF) kernel to transform the data, and then use Logistic Regression for classification.

## **Q13]**

Adaboost (Adaptive Boosting) is a boosting algorithm that fits a sequence of weak learners on repeatedly modified versions of the data. The predictions from all of them are then combined through a weighted majority vote (or sum) to produce the final prediction. It's worth noting that Adaboost is sensitive to noisy data and outliers.

Gradient Boosting, on the other hand, builds an additive model in a forward stage-wise fashion, allowing for the optimization of arbitrary differentiable loss functions. It generalizes the boosting process by allowing for the optimization of any differentiable loss function, rather than being limited to the exponential

loss function used in Adaboost. Gradient Boosting is also less sensitive to noisy data and outliers compared to Adaboost.

#### **Q14]**

In machine learning, the bias-variance trade-off is a fundamental concept that describes the relationship between the complexity of a model and its ability to generalize to new data.

Bias refers to the error that is introduced by approximating a real-world problem with a simplified model. A high bias model is oversimplified and may not capture the underlying patterns in the data, leading to poor performance on both the training and testing data.

The bias-variance trade-off is the balance between these two sources of error. As the complexity of a model increases, its variance typically increases and its bias decreases. Conversely, as the complexity of a model decreases, its variance decreases and its bias increases.

The goal of machine learning is to find the optimal balance between bias and variance that results in the lowest possible error on new, unseen data. This is often achieved through techniques such as regularization, cross-validation, and ensemble methods.

Regularization adds a penalty term to the model's objective function to discourage overfitting, while cross-validation involves splitting the data into multiple subsets and training the model on each subset to estimate its generalization performance. Ensemble methods, such as bagging and boosting, combine multiple models to reduce variance and improve generalization performance.

the bias-variance trade-off is a fundamental concept in machine learning that describes the relationship between the complexity of a model and its ability to generalize to new data. Finding the optimal balance between bias and variance is critical for building accurate and robust machine learning models.

#### **Q15]**

1. Linear Kernel: The linear kernel is the simplest kernel function used in SVM. It calculates the inner product between two input vectors and adds a bias term. The linear kernel is used when the data is linearly separable, meaning that there exists a hyperplane that can separate the two classes with zero error. The linear kernel is defined as:

$$k(x, y) = x^T y + b$$

where  $x$  and  $y$  are input vectors, and  $b$  is a bias term.

2. Radial Basis Function (RBF) Kernel: The RBF kernel is a popular kernel function used in SVM for non-linearly separable data. It measures the similarity between two input vectors based on the Euclidean distance between them. The RBF kernel is defined as:

$$k(x, y) = \exp(-\gamma \|x - y\|^2)$$

where  $x$  and  $y$  are input vectors,  $\gamma$  is a free parameter that controls the width of the Gaussian function, and  $\|x - y\|^2$  is the squared Euclidean distance between the two vectors.

3. Polynomial Kernel: The polynomial kernel is another kernel function used in SVM for non-linearly separable data. It calculates the inner product between two input vectors raised to a power, along with a bias term. The polynomial kernel is defined as:

$$k(x, y) = (\gamma x^T y + \text{coef0})^{\text{degree}}$$

where  $x$  and  $y$  are input vectors,  $\gamma$  is a free parameter that controls the influence of each feature,  $\text{coef0}$  is a bias term, and  $\text{degree}$  is the degree of the polynomial.

the linear kernel is used for linearly separable data, the RBF kernel is used for non-linearly separable data with a Gaussian-like similarity measure, and the polynomial kernel is used for non-linearly separable data with a polynomial-like similarity measure. The choice of kernel function depends on the nature of the data and the problem at hand.



## **STATISTICS WORKSHEET ANSWERS :**

**Q1]**

**Ans:** d) Expected

**Q2]**

**Ans:** c) Frequencies

**Q3]**

**Ans:** c) 6

**Q4]**

**Ans:** b) Chisquared distribution

**Q5]**

**Ans:** c) F Distribution

**Q6]**

**Ans:** b) Hypothesis

**Q7]**

**Ans:** a) Null Hypothesis

**Q8]**

**Ans:** a) Two tailed

**Q9]**

**Ans:** b) Research Hypothesis

**Q10]**

**Ans:** a) np