

## **WORKSHEET -5 MACHINE LEARNING**

1. R-squared is a commonly used measure of goodness of fit in regression. R-squared measures the proportion of variance in the dependent variable that is explained by the independent variables. R-squared values range from 0 to 1, with higher values indicating a better fit of the model to the data. On the other hand, Residual Sum of Squares (RSS) measures the sum of squared differences between the actual and predicted values of the dependent variable. Although it provides a measure of the lack of fit of the model, it is less interpretable than R-squared, as it does not provide a simple interpretation in terms of the proportion of variance explained by the model. Therefore, R-squared is generally considered to be a better measure of goodness of fit in regression than RSS.
2. In regression analysis, TSS (Total Sum of Squares) measures the total variance in the dependent variable. It is calculated as the sum of squared differences between the mean of the dependent variable and each observed value of the dependent variable. ESS (Explained Sum of Squares) measures the variance in the dependent variable that is explained by the independent variables. It is calculated as the sum of squared differences between the predicted values of the dependent variable (based on the regression model) and the mean of the dependent variable. RSS (Residual Sum of Squares) measures the variance in the dependent variable that is not explained by the independent variables. It is calculated as the sum of squared differences between the observed values of the dependent variable and the predicted values of the dependent variable.

The three metrics are related to each other as follows:

$$\text{TSS} = \text{ESS} + \text{RSS}$$

where TSS is the total sum of squares, ESS is the explained sum of squares, and RSS is the residual sum of squares.

3. Regularization refers to techniques that are used to calibrate machine learning models in order to minimize the adjusted loss function and prevent overfitting or underfitting. Using Regularization, we can fit our machine

learning model appropriately on a given test set and hence reduce the errors in it. In the context of machine learning, regularization is the process which regularizes or shrinks the coefficients towards zero. In simple words, regularization discourages learning a more complex or flexible model, to prevent overfitting. There are two main types of regularization techniques: Ridge Regularization and Lasso Regularization.

4. Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset. In other words, it measures the degree of mixture of elements with different class labels in a set. It is used as a splitting criterion in decision tree algorithms to determine the best split of a node in a decision tree. The Gini impurity is calculated as the probability of misclassifying a randomly chosen element if the labels were randomly distributed among the elements. The Gini impurity ranges from 0 (indicating a pure set) to 1 (indicating an impure set with a maximum degree of mixture).
5. Yes, unregularized decision trees are prone to overfitting. Overfitting occurs when a model is too complex and captures the noise in the training data, rather than the underlying pattern. This leads to poor performance on unseen data, as the model is not general enough to make accurate predictions. Unregularized decision trees are particularly prone to overfitting because they continue to split the data until each leaf node contains only one class label or a minimum number of samples. This results in highly complex trees with many levels and branches that are very sensitive to the training data. The model can memorize the training data, leading to overfitting. Therefore, it is often necessary to regularize decision trees, such as by setting a maximum depth or minimum number of samples required to split a node, in order to avoid overfitting and improve the generalization performance of the model.
6. Ensemble techniques in machine learning are methods that combine the predictions of multiple models to improve the overall accuracy and stability of the predictions. The idea behind ensemble techniques is to train multiple models with different architectures or learning algorithms, and then aggregate their predictions in some way to form a final prediction. There are several popular ensemble techniques, including: Bagging (Bootstrap Aggregating): This involves training multiple models on bootstrapped samples of the training data and aggregating their predictions through

majority voting or averaging. Boosting: This involves training a sequence of models, with each model focusing on the samples that were misclassified by the previous models. The final prediction is obtained by weighting the predictions of each model. Random Forest: This is an extension of bagging that involves training decision trees on bootstrapped samples of the training data, and aggregating their predictions through majority voting. Stacking: This involves training multiple models and using their predictions as input features to train a meta-model, which makes the final prediction. Ensemble techniques are widely used in practice due to their ability to provide improved accuracy and stability compared to individual models, especially when the individual models have high variance or low bias.

7. Bagging (Bootstrap Aggregating) and Boosting are two popular ensemble techniques in machine learning that combine the predictions of multiple models to improve the overall accuracy and stability of the predictions. However, the two techniques have different approaches to combining the models. Bagging: Bagging involves training multiple models on bootstrapped samples of the training data and aggregating their predictions through majority voting or averaging. The goal of bagging is to reduce the variance of the model by aggregating the predictions of multiple models that are trained on different samples of the data. The individual models in a bagging ensemble are typically of the same type and trained on different samples of the training data, leading to diverse predictions that can be aggregated to reduce the variance. Boosting: Boosting involves training a sequence of models, with each model focusing on the samples that were misclassified by the previous models. The final prediction is obtained by weighting the predictions of each model. The goal of boosting is to reduce the bias of the model by focusing the training of each model on the samples that are difficult to classify. The individual models in a boosting ensemble are typically of the same type and trained on the same data, but with different weights assigned to each sample based on the performance of the previous models. In summary, bagging aims to reduce the variance of the model, while boosting aims to reduce the bias of the model. Both techniques are used to improve the stability and accuracy of the predictions, and are often used in combination to achieve the best results.
8. In random forests, an out-of-bag (OOB) error is a measure of the accuracy of the model that is estimated without using the test data. The idea behind OOB error is that in a random forest, each tree is trained on a bootstrapped sample of the training data, meaning that some samples are not used in the training of a particular tree. These samples that are not used in the training of a particular tree can be used to estimate the accuracy of that tree without using

any test data. The OOB error is calculated by first making predictions for the samples that were not used in the training of a particular tree, and then comparing these predictions with the actual values. The average of the OOB errors across all trees in the random forest gives an estimate of the overall accuracy of the model. The OOB error can be used for model selection and for estimating the generalization error of the random forest model. It is also used for variable importance measures, as the difference in OOB error between the full model and a model where a feature is permuted can be used to estimate the importance of that feature.

9. K-fold cross-validation is a technique used to evaluate the performance of a machine learning model by dividing the dataset into K parts, or folds, of roughly equal size. The model is trained on K-1 of the folds and tested on the remaining one. This procedure is repeated K times, with each fold used exactly once as the test set. The results of the K test sets are then averaged to give an overall estimate of the model's performance. K-fold cross-validation provides a more robust estimate of the model's performance compared to using a single train-test split, as it averages the performance over multiple different train-test splits. The choice of the number of folds (K) affects the trade-off between the bias and the variance of the estimate. A lower value of K means that the test set is smaller, which results in a high variance estimate, while a larger value of K results in a high bias estimate. A common choice for K is 10 or 5. K-fold cross-validation is a widely used technique for model selection and for estimating the generalization error of a model, as it provides a more reliable estimate of the model's performance on new, unseen data. It can also be used to tune the hyperparameters of the model, by selecting the hyperparameters that result in the best average performance across the K folds.
10. Hyperparameter tuning is the process of optimizing the hyperparameters of a machine learning model to obtain the best performance on a given dataset. Hyperparameters are parameters that are set prior to training a model, as opposed to model parameters, which are learned from the data during training. Examples of hyperparameters include the learning rate for gradient descent, the number of trees in a random forest, or the regularization strength in a linear regression model. Hyperparameter tuning is important because the performance of a machine learning model is often highly sensitive to the values of its hyperparameters. For example, a model with too large a learning rate may converge slowly or not at all, while a model with too small a learning rate may converge to a suboptimal solution. Similarly, a model with too many trees may overfit the data, while a model with too few trees may underfit the data. The goal of hyperparameter tuning is to find the

hyperparameters that result in the best performance of the model on a given dataset. This is usually done by searching over a range of hyperparameters, training the model with each combination of hyperparameters, and evaluating its performance on a validation set. The hyperparameters that result in the best performance are then selected for use on the test set. Hyperparameter tuning can be a time-consuming and computationally expensive process, but it is an important step in obtaining the best possible performance from a machine learning model. Some techniques, such as grid search and random search, can be used to efficiently search over the hyperparameter space, while others, such as Bayesian optimization, use more sophisticated methods to find the best hyperparameters.

11. A large learning rate in gradient descent can result in several issues:

Oscillation: If the learning rate is too large, the gradient update step can be too big, causing the model to overshoot the minimum and oscillate between high and low values, never settling on a good solution. Divergence: A learning rate that is too large can cause the cost function to increase instead of decrease, resulting in the model never converging to a solution. Slow convergence: If the learning rate is too small, gradient descent will converge very slowly, as the model takes small steps that require many iterations to reach the minimum. However, if the learning rate is too large, the model may converge too quickly to a suboptimal solution. Instability: A large learning rate can cause the model to be unstable, with the parameters changing erratically and the model never settling on a good solution. Overfitting: A large learning rate can cause the model to overfit the training data, as it quickly tries to fit the data and ends up memorizing the noise instead of generalizing to the underlying patterns. Therefore, it's important to choose a learning rate that is not too small or too large. A common approach is to start with a large learning rate and gradually decrease it over time. Another common approach is to use a learning rate schedule, where the learning rate is adjusted automatically as the training progresses, such as with step decay or adaptive learning rates.

12. Logistic Regression is a linear classifier, meaning it assumes a linear relationship between the input features and the target variable. It models the probability of a target variable belonging to a certain class as a linear combination of the input features. For non-linear data, the relationship between the input features and target variable is not a straight line, and Logistic Regression may not accurately capture this relationship. In such cases, other non-linear classifiers like decision trees, random forests, or

support vector machines (SVMs) with non-linear kernels may perform better. However, it is possible to transform the input features in such a way that the relationship between the input features and target variable becomes linear. For example, adding polynomial features to the input data, or using logarithmic transformations, can help make the data linearly separable, allowing Logistic Regression to be applied.

13. Adaboost and Gradient Boosting are both ensemble techniques used for improving the performance of weak models, but there are some differences between the two:  
**Algorithm:** Adaboost is an adaptive boosting algorithm, where the weights of the misclassified examples are adjusted in each iteration to give more importance to the difficult examples. Gradient Boosting is a gradient descent optimization-based algorithm, where the objective is to minimize the loss function by adding weak models in a stagewise manner.  
**Loss Function:** Adaboost uses a classification loss function, such as the exponential loss, which measures the confidence in the prediction. Gradient Boosting uses a differentiable loss function, such as mean squared error or mean absolute error, which can be optimized using gradient descent.  
**Model Complexity:** Adaboost uses a set of simple models, such as decision stumps, which have a single split. Gradient Boosting allows for more complex models, such as decision trees, which have multiple splits.  
**Computational Complexity:** Adaboost is computationally simpler than Gradient Boosting, as it requires only a simple linear optimization problem to be solved at each iteration. Gradient Boosting requires solving a non-linear optimization problem, which can be more computationally intensive.  
In summary, Adaboost is a simple and fast algorithm that is suitable for binary classification problems, while Gradient Boosting is a more flexible and powerful algorithm that can be applied to regression and classification problems.
14. The bias-variance tradeoff is a fundamental concept in machine learning that refers to the tradeoff between the accuracy of a model's predictions and its complexity. Bias refers to the error that is introduced by approximating a real-world problem with a simpler model. Models with high bias tend to underfit the data, meaning they don't capture the true underlying relationship between the input features and the target variable. High bias models have high error on both the training and testing data. Variance refers to the error that is introduced by a model that is too complex and fits the training data too closely. Models with high variance tend to overfit the data, meaning they capture noise or random fluctuations in the training data. High variance models have low error on the training data but high error on the testing data. The bias-variance tradeoff means that increasing the complexity of a model

reduces bias but increases variance, and vice versa. The goal in machine learning is to find a balance between bias and variance to achieve a model that makes accurate predictions and generalizes well to unseen data. This tradeoff can be managed by techniques such as regularization, early stopping, and ensembling.

15. Linear kernel: The linear kernel is the simplest kernel function and works well when the data is linearly separable. It simply computes the dot product between two vectors, representing the feature values of two instances. The linear kernel is used in linear Support Vector Machines (SVMs). Radial basis function (RBF) kernel: The RBF kernel is a non-linear kernel function that is used when the data is not linearly separable. The RBF kernel maps the input features into a higher-dimensional feature space through a radial basis function, which helps to separate the classes in this transformed space. Polynomial kernel: The polynomial kernel is another non-linear kernel function that maps the input features into a higher-dimensional feature space through a polynomial transformation. It can be used for polynomial classification problems. The polynomial kernel is defined as the dot product of two instances raised to a power, which is a parameter of the kernel. The degree of the polynomial determines the degree of the polynomial transformation.