

## MACHINE LEARNING ASSIGNMENT - 5

Q1 to Q15 are subjective answer type questions, Answer them briefly.

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:- Both R-squared and Residual Sum of Squares (RSS) are commonly used measures to evaluate the goodness of fit of a regression model, but they capture different aspects of the model's performance.

R-squared is a statistical measure that indicates the proportion of variance in the dependent variable that is explained by the independent variables in the model. It ranges from 0 to 1, with higher values indicating a better fit of the model to the data. R-squared is often used as a measure of the overall quality of the regression model because it provides an indication of how well the model fits the data.

On the other hand, RSS measures the sum of the squared differences between the predicted values and the actual values in the dataset. It is a measure of the amount of unexplained variation in the data that the model cannot account for. Lower values of RSS indicate a better fit of the model to the data.

In general, R-squared is considered a better measure of goodness of fit than RSS because it takes into account the proportion of variance in the dependent variable that is explained by the independent variables. However, R-squared has some limitations, as it can be sensitive to the number of variables in the model, and it may not always provide a complete picture of the model's performance.

Therefore, it is recommended to use both R-squared and RSS to evaluate the goodness of fit of a regression model. R-squared can provide an overall assessment of the model's performance, while RSS can help identify specific areas where the model is not fitting the data well.

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:- In regression, TSS (Total Sum of Squares), ESS (Explained Sum of Squares), and RSS (Residual Sum of Squares) are three important measures that are used to evaluate the goodness of fit of a regression model.

TSS (Total Sum of Squares) is the total variation in the dependent variable that is captured by the model. It represents the sum of the squared differences between the actual values of the dependent variable and the mean of the dependent variable.

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

RSS (Residual Sum of Squares) is the variation in the dependent variable that is not explained by the independent variable(s) in the model. It represents the sum of the squared differences between the actual values of the dependent variable and the predicted values of the dependent variable.

The equation relating these three metrics is:

$$TSS = ESS + RSS$$

This equation represents the partitioning of the total variation in the dependent variable into the variation that is explained by the independent variable(s) in the model (ESS) and the variation that is not explained by the independent variable(s) in the model (RSS).

In other words, TSS represents the total variation in the dependent variable, ESS represents the variation that is explained by the model, and RSS represents the variation that is not explained by the model. By comparing the relative magnitudes of ESS and RSS, we can assess the goodness of fit of the regression model.

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

Answer:- Regularization is a technique used in machine learning to prevent overfitting of the model. Overfitting occurs when the model is too complex and captures noise in the data, resulting in poor generalization to new data. Regularization helps to simplify the model by adding a penalty term to the cost function, which encourages the model to have smaller weights and biases.

The need for regularization arises because, in many real-world scenarios, the amount of available data is limited, and the complexity of the model can quickly exceed the amount of data available. In such cases, a model with high variance can be susceptible to overfitting the training data, which results in poor performance on new, unseen data.

Regularization techniques, such as L1 regularization (Lasso) and L2 regularization (Ridge), can help to prevent overfitting by adding a penalty term to the cost function. These techniques add a constraint to the optimization problem that forces the weights to be small, which reduces the complexity of the model. This results in a model that is more likely to generalize well to new data.

Regularization can also help to improve the stability and interpretability of the model, as it reduces the sensitivity of the model to small changes in the data, and can also help to eliminate irrelevant features from the model.

In summary, the need for regularization in machine learning arises from the desire to prevent overfitting, improve generalization to new data, increase model stability and interpretability, and reduce model complexity.

### 4. What is Gini-impurity index?

Answer:- The Gini impurity index is a measure of the impurity or the heterogeneity of a set of samples. It is commonly used as a criterion to build decision trees in classification problems.

The Gini impurity index measures the probability of misclassifying a randomly chosen sample from a set if it were randomly labeled according to the distribution of labels in the set. In other words, it calculates the probability of two randomly chosen elements from the set being incorrectly labeled.

The formula for the Gini impurity index is:

$$\text{Gini} = 1 - \sum p^2$$

where  $p$  is the proportion of samples in the set that belong to a particular class or category.

The Gini impurity index ranges from 0 to 1, with 0 representing a completely pure set, where all samples belong to the same class, and 1 representing a completely impure set, where the samples are evenly distributed among all classes.

When building decision trees, the Gini impurity index is used to determine the best split at each node of the tree. The split that results in the lowest Gini impurity index is chosen as the best split.

In summary, the Gini impurity index is a measure of the impurity or heterogeneity of a set of samples and is commonly used as a criterion to build decision trees in classification problems.

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

Answer:- Yes, unregularized decision trees can be prone to overfitting. Overfitting occurs when the model is too complex and fits the training data too well, capturing noise and idiosyncrasies in the data rather than general patterns.

Decision trees are prone to overfitting because they have high variance and can create complex decision boundaries that may fit the training data perfectly but do not generalize well to new, unseen data. As the tree grows deeper and more complex, it becomes increasingly likely to overfit the training data.

To mitigate overfitting, regularization techniques such as pruning, limiting the depth of the tree, and reducing the minimum number of samples required to split a node can be used. Additionally, ensemble methods such as random forests and boosting can be used to combine multiple decision trees and reduce the risk of overfitting.

## 6. What is an ensemble technique in machine learning?

Answer:- Ensemble techniques in machine learning are methods that combine multiple models to improve the overall performance of the prediction. The idea behind ensemble methods is that by combining several models that individually may not perform well, a better and more accurate prediction can be made.

There are several types of ensemble techniques, including:

**Bagging (Bootstrap Aggregating):** Bagging involves training multiple instances of the same model on different subsets of the training data, then aggregating their predictions by taking the average or majority vote.

**Boosting:** Boosting involves training a sequence of models, with each subsequent model focusing on the training instances that were previously misclassified by the previous models.

**Stacking:** Stacking involves training multiple models, then using the predictions of these models as input features to a higher-level model.

**Random forests:** Random forests are an extension of bagging that build a large number of decision trees and aggregate their predictions.

Ensemble methods are particularly useful when individual models have high variance or when there is no single model that can accurately capture the complexity of the data. By combining multiple models, ensemble methods can reduce the risk of overfitting and improve the accuracy and robustness of the prediction.

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

Answer:- Bagging (Bootstrap Aggregating) and Boosting are two popular ensemble techniques used in machine learning. While they both involve

combining multiple models to improve the overall performance, there are some key differences between the two methods:

**Training Process:** In Bagging, multiple instances of the same model are trained on different subsets of the training data. Each model is trained independently of the others, and the predictions are aggregated by taking the average or majority vote. In Boosting, a sequence of models is trained, with each subsequent model focusing on the training instances that were previously misclassified by the previous models.

**Model Diversity:** In Bagging, the focus is on creating diversity among the models by training them on different subsets of the training data. In Boosting, the focus is on improving the accuracy of the prediction by training models sequentially, with each model building on the weaknesses of the previous models.

**Handling of Outliers:** Bagging is more robust to outliers since the models are trained on different subsets of the data. Boosting is more sensitive to outliers since it focuses on correcting the misclassified instances from previous models.

**Risk of Overfitting:** Bagging typically has a lower risk of overfitting since it involves averaging the predictions of multiple models. Boosting, on the other hand, has a higher risk of overfitting since it involves training models sequentially and potentially memorizing the training data.

In summary, Bagging and Boosting are both useful techniques for improving the performance of machine learning models. Bagging focuses on creating diversity among models and reducing the risk of overfitting, while Boosting focuses on improving accuracy by correcting the misclassifications of previous models.

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

Answer:- In random forests, the out-of-bag (OOB) error is a metric used to estimate the performance of the model on unseen data.

When constructing a random forest model, each decision tree is trained on a random subset of the training data, leaving a portion of the data unused, known as the out-of-bag samples. The OOB error is the average error of each decision tree on their respective out-of-bag samples.

The OOB error is a useful metric for assessing the performance of the model without the need for a separate validation set. It provides an unbiased estimate of the model's generalization error and can be used to tune hyperparameters, such as the number of trees in the forest or the maximum depth of the trees.

Furthermore, the OOB error can be used to perform feature selection, as it indicates the importance of each feature in the model. The importance of a feature is measured by the reduction in the OOB error when the feature is randomly permuted, disrupting the relationship between the feature and the response variable.

Overall, the OOB error is a valuable tool in random forests for assessing the model's performance, selecting hyperparameters, and performing feature selection.

#### 9. What is K-fold cross-validation?

Answer:- K-fold cross-validation is a popular technique used in machine learning to evaluate the performance of a model and to tune its hyperparameters. It involves splitting the data into K equally sized folds, where K is a user-specified parameter.

The algorithm then trains the model K times, using K-1 folds as the training set and the remaining fold as the validation set. In each iteration, a different fold is used as the validation set, and the other K-1 folds are used for training. This

process is repeated K times, with each fold serving as the validation set exactly once.

The performance of the model is evaluated by taking the average of the performance metrics obtained in each fold. This method helps to reduce the variance of the performance estimate, as each observation is used both for training and validation, and it helps to ensure that the model can generalize well to new data.

K-fold cross-validation can also be used to tune the hyperparameters of the model, such as the regularization parameter or the learning rate, by evaluating the model's performance on different hyperparameter values across different folds. By selecting the hyperparameters that result in the best average performance across the K folds, one can obtain a model that is more likely to generalize well to unseen data.

Overall, K-fold cross-validation is a useful and widely used technique in machine learning for evaluating model performance, selecting hyperparameters, and assessing the model's ability to generalize to new data.

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

Answer:- Hyperparameter tuning is the process of finding the optimal hyperparameters for a machine learning algorithm. Hyperparameters are parameters that are not learned during the training process, such as the learning rate, regularization parameter, number of hidden layers in a neural network, or the number of trees in a random forest.

Hyperparameter tuning is done to find the best combination of hyperparameters that can optimize the performance of a machine learning algorithm. The goal is to find hyperparameters that can maximize the performance metrics, such as accuracy, precision, recall, or F1-score, on a validation set.

The process of hyperparameter tuning involves the following steps:



Define a range of hyperparameters: A range of hyperparameters is defined, and a search algorithm is used to explore the space of possible values for these hyperparameters.

Select a search algorithm: There are different search algorithms that can be used for hyperparameter tuning, such as grid search, random search, or Bayesian optimization.

Split the data into training and validation sets: The data is split into a training set and a validation set, and the model is trained on the training set using different hyperparameters.

Evaluate the model performance: The performance of the model is evaluated on the validation set using different performance metrics, such as accuracy, precision, recall, or F1-score.

Select the best hyperparameters: The hyperparameters that result in the best performance on the validation set are selected, and the final model is trained on the entire training set using these hyperparameters.

Overall, hyperparameter tuning is an essential step in machine learning as it helps to improve the performance of the model and ensures that the model can generalize well to new data. It can be a time-consuming process, but it can lead to significant improvements in the performance of the model.

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

Answer:- If the learning rate is too large in gradient descent, it can lead to several issues, including:

Divergence: A large learning rate can cause the algorithm to overshoot the minimum point and bounce back and forth between the two sides of the

minimum point. This can cause the algorithm to diverge and fail to converge to a solution.

Slow convergence: If the learning rate is too large, the algorithm may take a long time to converge, as the oscillations around the minimum point slow down the convergence.

Overshooting the minimum point: A large learning rate can cause the algorithm to overshoot the minimum point and move away from the optimal solution. This can result in a less accurate and suboptimal model.

Instability: A large learning rate can lead to instability in the optimization process, making it difficult to find the optimal solution.

To avoid these issues, it is important to choose an appropriate learning rate for the problem at hand. This can be done by performing a grid search or using an adaptive learning rate algorithm, such as AdaGrad, RMSProp, or Adam, which adjust the learning rate automatically during the training process.

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

Answer:- Logistic regression is a linear model that is commonly used for binary classification problems, where the response variable is binary, such as 0 or 1. Logistic regression assumes that the decision boundary between the two classes is a linear function of the input variables. This means that logistic regression may not perform well on datasets that have a non-linear decision boundary.

If the dataset has a non-linear decision boundary, logistic regression may not be able to capture the complex relationships between the input variables and the response variable. In such cases, it may be necessary to use a non-linear model, such as a support vector machine (SVM), decision tree, or a neural network, to model the non-linear relationships between the input variables and the response variable.

However, logistic regression can still be used for non-linear classification problems if the input variables are transformed using non-linear transformations, such as polynomial features, before fitting the model. This approach is known as polynomial logistic regression or generalized additive models (GAMs). Polynomial logistic regression can capture non-linear relationships between the input variables and the response variable by adding polynomial terms to the linear model. However, the complexity of the model increases with the number of polynomial terms, which can lead to overfitting if the model is not properly regularized.

Overall, while logistic regression can be used for non-linear classification problems using polynomial transformations, it may not be the best choice for highly non-linear datasets. Other non-linear models may perform better and provide more accurate predictions.

### 13. Differentiate between Adaboost and Gradient Boosting.

Answer:- Adaboost and Gradient Boosting are two popular ensemble learning techniques used in machine learning. While both methods aim to improve the performance of weak learners, they differ in their approach to achieve this goal. Here are the main differences between Adaboost and Gradient Boosting:

Training process: Adaboost uses a weighted combination of weak learners to improve the performance of the model, whereas Gradient Boosting trains a sequence of weak learners that are combined to form a strong learner.

Loss function: Adaboost focuses on minimizing the exponential loss function, whereas Gradient Boosting minimizes any differentiable loss function, such as squared loss or absolute loss.

Weight updates: In Adaboost, the weights of the training samples are updated after each iteration based on their misclassification rate, whereas in Gradient Boosting, the residuals of the previous iteration are used to update the weights of the training samples.

Learning rate: Adaboost uses a fixed learning rate for all iterations, whereas Gradient Boosting uses a smaller learning rate for each iteration, which helps to prevent overfitting.

Handling outliers: Adaboost is sensitive to outliers, which can lead to overfitting, whereas Gradient Boosting can handle outliers and is less prone to overfitting.

Overall, both Adaboost and Gradient Boosting are powerful ensemble learning techniques that can significantly improve the performance of machine learning models. The choice between the two methods depends on the specific characteristics of the problem at hand, such as the nature of the data and the type of loss function that is being minimized.

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

Answer:- Bias-variance tradeoff is a fundamental concept in machine learning that describes the tradeoff between the model's ability to fit the training data (bias) and its ability to generalize to new, unseen data (variance).

Bias refers to the error that arises from the simplifying assumptions made by a model to make the target function easier to learn. Models with high bias are often too simple and unable to capture the complex patterns present in the data. These models typically underfit the data, leading to poor training and test performance.

Variance, on the other hand, refers to the error that arises from the model's sensitivity to small fluctuations in the training data. Models with high variance are often too complex and overfit the training data, leading to poor test performance.

The goal of machine learning is to find a balance between bias and variance, such that the model can generalize well to new data while still fitting the

training data. This balance can be achieved by tuning the complexity of the model or by applying regularization techniques.

In general, complex models such as deep neural networks have high variance and low bias, while simpler models such as linear regression have high bias and low variance. Regularization techniques such as L1 and L2 regularization can help to reduce the variance of complex models, while increasing the bias of simpler models can help to reduce their variance.

Finding the right balance between bias and variance is an important part of building a successful machine learning model. A model with too much bias will have poor accuracy, while a model with too much variance will overfit the training data and have poor generalization performance.

15. Give short description each of Linear, RBF, Polynomial kernels used in SVM?

Answer:- Support Vector Machines (SVMs) are a type of supervised learning algorithm used for classification and regression tasks. SVMs use kernels to transform the input data into a higher-dimensional space, where it is easier to find a linear separation boundary between classes.

Here are the three commonly used types of kernels in SVM:

Linear kernel: The linear kernel is the simplest kernel and is used when the data is linearly separable. It works by computing the dot product between the input vectors and is given by  $K(x, y) = x.T * y$ .

RBF (Radial Basis Function) kernel: The RBF kernel is a non-linear kernel that works well when the data is not linearly separable. It maps the input data into a high-dimensional space and computes the similarity between two data points using the Euclidean distance. The RBF kernel is given by  $K(x, y) = \exp(-\gamma * ||x-y||^2)$ , where  $\gamma$  is a hyperparameter that controls the width of the Gaussian kernel.

Polynomial kernel: The polynomial kernel is also a non-linear kernel that works well when the data is not linearly separable. It maps the input data into a higher-dimensional space using a polynomial function. The polynomial kernel is given by  $K(x, y) = (x.T * y + r)^d$ , where  $r$  and  $d$  are hyperparameters that control the degree of the polynomial.

Overall, SVMs are powerful machine learning algorithms that can be used to solve a wide range of classification and regression problems. The choice of kernel depends on the nature of the data and the problem at hand.