

# MACHINE LEARNING

## ASSIGNMENT - 2

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?

Ans 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 model performance.

### 1. R-squared (Coefficient of Determination):

- R-squared measures the proportion of the variance in the dependent variable that is explained by the independent variables in the model.
- It ranges from 0 to 1, where 0 indicates that the model does not explain any variability in the dependent variable, and 1 indicates that the model explains all the variability.
- Higher values of R-squared indicate a better fit of the model to the data.

### 2. Residual Sum of Squares (RSS):

- RSS measures the total squared difference between the observed values and the predicted values by the regression model.
- It represents the unexplained variance in the dependent variable by the independent variables in the model.
- Lower values of RSS indicate a better fit of the model to the data.

### Which one is better?

- R-squared is often preferred for comparing models or assessing the overall explanatory power of the model because it provides a standardized measure of goodness of fit that is easy to interpret.
- However, RSS can also be useful, especially if you're interested in understanding the magnitude of the errors or residuals in the model.

- In practice, it's often a good idea to consider both measures together. A high R-squared coupled with a low RSS indicates that the model explains a large proportion of the variance in the dependent variable while also minimizing the discrepancy between observed and predicted values.

In summary, while R-squared is more commonly reported and interpreted, both R-squared and RSS provide valuable insights into the performance of a regression model, and using them together can offer a more comprehensive evaluation.

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 analysis TSS Total sum of squares ESS explained sum of square RSS residual sum of square is an important metric to evaluate the Goodness of fit in regression model

**TSS Total Sum of Squares**

TSS measures the total variability in dependent variable (y) before the regression model is applied. It represents the total deviation of the observed values from the mean of the dependent variable

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2$$

Where  $y_i$  represents the observed values of the dependent variable,  $\bar{y}$  is the mean of the dependent variable and  $n$  is the number of observations.

**ESS Explained Sum of Squares**

ESS measures the variability in the dependent variable that is explained by the regression model. It quantifies how much of the total variability in the dependent variable is accounted for by the independent variables in the regression model.

Mathematically, it is calculated as

$$ESS = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$\hat{y}_i$  represents the predicted values of the dependent variable obtained from the regression model, and  $\bar{y}$  is the mean of the dependent variable.

**RSS Residual sum of squares**

RSS measures the unexplained variability in the dependent variable after the regression model has been applied. It quantifies the difference between the observed values of the dependent variable and the value predicted by the regression model. Mathematically, it is calculated as

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

where  $y_i$  represents the observed values of the dependent variable,  $\hat{y}_i$  represents the predicted values of the dependent variable obtained from the regression model, and  $n$  is the number of observations.

The relationship between these three metrics is given by:

$$TSS = ESS + RSS$$

In other word, the total variability in the dependent variable (TSS) is equal to the sum of the variability explained by the regression model (ESS) and the unexplained variability(RSS).

### 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 ability of models. Overfitting occurs when a model learns to fit the training data too closely, capturing noise or random fluctuations that are not representative of the true underlying relationship between the features and the target variable. Regularization helps address this issue by adding a penalty term to the model's objective function,

which discourages overly complex models with large parameter values.

The need for regularization arises due to several reasons:

1. **Prevention of Overfitting:** Regularization helps prevent overfitting by penalizing overly complex models. By constraining the model's complexity, regularization encourages it to capture the underlying patterns in the data rather than memorizing noise or outliers in the training set.
2. **Improvement of Generalization:** Regularization improves the generalization ability of models, allowing them to perform well on unseen data. By reducing overfitting, regularization helps ensure that the model's performance is more representative of its true predictive ability on new, unseen examples.
3. **Handling Multicollinearity:** In regression models with highly correlated features (multicollinearity), regularization techniques such as Ridge Regression or Lasso Regression can help stabilize parameter estimates and improve model interpretability.
4. **Dealing with High-Dimensional Data:** Regularization is particularly useful when working with high-dimensional data, where the number of features is much larger than the number of samples. In such cases, regularization can help select relevant features and avoid overfitting due to the curse of dimensionality.
5. **Robustness to Outliers:** Regularization can make models more robust to outliers by reducing their influence on parameter estimation. Techniques like Ridge Regression and Lasso Regression shrink the parameter estimates towards zero, reducing the impact of outliers on the final model.

Overall, regularization is a fundamental technique in machine learning for improving the performance and robustness of models, especially in scenarios where overfitting is a concern or when working with high-dimensional data.

4. What is Gini-impurity index?

Answer Gini-impurity index is a measure used in Decision tree algorithms particularly in Classification tasks to evaluate how well a split separates the classes within the data. It quantifies the impurity of a set of data points, where lower value indicate pure nodes with predominantly one class higher value indicate impurer nodes with multiple classes

$$Gini(p) = 1 - \sum_{i=1}^k p_i^2$$

Where:

- $p$  is a vector containing the proportions of each class in the set.
- $k$  is the number of classes.
- $p_i$  is the proportion of data points belonging to class  $i$  in the set.

The Gini impurity index ranges from 0 to 1. A Gini impurity of 0 indicates that the set contains only one class (perfectly pure), while a Gini impurity of 1 indicates an equal distribution of all classes (maximum impurity).

In the context of decision trees, the algorithm iterates over different features and split points to find the split that minimizes the weighted sum of Gini impurity for the resulting child nodes. This process continues recursively until a stopping criterion is met, such as reaching a maximum depth or minimum number of data points in a node.

The Gini impurity index is commonly used in decision tree algorithms like CART (Classification and Regression Trees) to determine the optimal splits for building the tree. It provides a criterion for selecting the best features and split points that lead to more homogeneous subsets, ultimately leading to a more accurate classification model.

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

Answer

Yes, unregularized decision trees are prone to overfitting.

Decision trees have a tendency to create complex models that fit the training data extremely well but may fail to generalize to unseen data. This overfitting occurs due to several reasons:

1. **High Variance:** Decision trees are capable of capturing intricate patterns in the training data, including noise and outliers. Without any constraints on the tree's structure, it can grow excessively deep to fit the training data perfectly, leading to high variance. As a result, small variations or noise in the training data can result in different tree structures, making the model sensitive to the specific training dataset.
2. **Memorization of Training Data:** Decision trees can potentially memorize the training data by creating leaf nodes corresponding to each training instance, especially in the absence of regularization. This memorization of training data does not generalize well to new data and leads to poor performance on unseen data.
3. **Lack of Pruning:** Unregularized decision trees may not undergo pruning, a process of reducing the size of the tree by removing nodes that provide little predictive power. Without pruning, the tree may grow excessively large and complex, capturing noise and irrelevant features in the training data.
4. **Sensitive to Data Distribution:** Decision trees can be sensitive to the distribution of the training data. They may split the data based on features that have no predictive power, especially if those features happen to separate the training data well but do not generalize to new data.

To mitigate overfitting in decision trees, various regularization techniques can be applied, such as:

- Limiting the maximum depth of the tree.
- Setting a minimum number of samples required to split a node.
- Pruning the tree after it has been built to remove unnecessary nodes.

- Using ensemble methods like Random Forests or Gradient Boosted Trees, which combine multiple decision trees to reduce overfitting.

By applying these regularization techniques, decision trees can generalize better to new data and avoid overfitting while still maintaining their interpretability and ease of use.

6.What is an ensemble technique in machine learning?

Answer

In machine learning, an ensemble technique is a method that combines the predictions of multiple individual models to produce a single, more robust prediction. Ensemble methods often outperform single models by leveraging the diversity of the individual models and combining their predictions in a way that reduces bias and variance, leading to improved generalization performance.

There are several types of ensemble techniques, including:

1. **Bagging (Bootstrap Aggregating)**: Bagging involves training multiple instances of the same base learning algorithm on different subsets of the training data. Each model is trained independently, typically using random sampling with replacement (bootstrap sampling). The final prediction is obtained by averaging (for regression) or taking a majority vote (for classification) of the predictions made by individual models. Random Forest is a popular example of a bagging ensemble technique, which uses decision trees as base learners.
2. **Boosting**: Boosting is an iterative ensemble technique that trains a sequence of weak learners sequentially, with each learner focusing on the instances that were misclassified by the previous ones. In boosting, each subsequent model in the sequence is trained to correct the errors made by the previous models, leading to a stronger ensemble model. Gradient Boosting Machines (GBM) and AdaBoost are common boosting algorithms.
3. **Stacking (Stacked Generalization)**: Stacking combines the predictions of multiple base models by training a meta-learner (or

blender) on the outputs of these base models. Instead of simply averaging or voting on the predictions, stacking learns to combine the strengths of individual models by training a higher-level model on their predictions. Stacking often leads to improved performance by capturing complementary patterns from diverse base models.

4. **Voting:** Voting combines the predictions of multiple models by taking a simple majority vote (for classification tasks) or averaging (for regression tasks). Different types of voting schemes include hard voting, where predictions are made based on the mode (most common class) of the individual predictions, and soft voting, where predictions are made based on the average probabilities predicted by the individual models.

Ensemble techniques are widely used in practice and have been shown to improve the robustness, accuracy, and stability of machine learning models across various domains and tasks. They are especially effective when individual models in the ensemble are diverse and make uncorrelated errors.

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

Answer

Bagging and boosting are both ensemble techniques used in machine learning, but they differ in their approach to combining multiple models and how they handle training data.

#### 1. **Bagging (Bootstrap Aggregating):**

- Bagging involves training multiple instances of the same base learning algorithm on different subsets of the training data. Each model is trained independently, typically using random sampling with replacement (bootstrap sampling).
- The final prediction is obtained by mean (for regression) or taking a majority vote (for classification) of the predictions made by individual models.
- Bagging helps reduce variance and overfitting by introducing diversity among the models. Since each model is trained on



a random subset of the data, they are less likely to overfit to the training set and more likely to generalize well to unseen data.

- Random Forest is a popular example of a bagging ensemble technique, which uses decision trees as base learners.

## 2. **Boosting:**

- Boosting is an iterative ensemble technique that trains a sequence of weak learners sequentially, with each learner focusing on the instances that were misclassified by the previous ones.
- In boosting, each subsequent model in the sequence is trained to correct the errors made by the previous models, leading to a stronger ensemble model.
- Boosting aims to reduce bias and improve accuracy by iteratively refining the model's predictions and focusing on difficult-to-classify instances.
- Common boosting algorithms include Gradient Boosting Machines (GBM) and AdaBoost.

In summary, the main differences between bagging and boosting are:

- Bagging trains multiple models independently and combines their predictions through averaging or voting, while boosting trains a sequence of models iteratively, with each model focusing on the mistakes of the previous ones.
- Bagging aims to reduce variance and overfitting by introducing diversity among models, while boosting aims to reduce bias and improve accuracy by iteratively refining the model's predictions

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

Answer In random forests, the out-of-bag (OOB) error is an estimation of the model's performance on unseen data without the need for a separate validation set. It is calculated using the data points that were not included in the bootstrap sample

(training set) used to train each individual tree in the random forest.

Here's how the out-of-bag error estimation works in random forests:

1. **Bootstrap Sampling:** Random forests use bootstrap sampling to create multiple training datasets. Each tree in the forest is trained on a different bootstrap sample, which is obtained by randomly selecting data points from the original dataset with replacement. Some data points may appear multiple times in the bootstrap sample, while others may not be selected at all.
2. **Out-of-Bag Data:** Since bootstrap sampling involves randomly selecting data points with replacement, approximately one-third of the original data points are left out of each bootstrap sample, on average. These left-out data points are called out-of-bag (OOB) data for a particular tree.
3. **OOB Error Calculation:** For each data point in the original dataset, the random forest calculates its prediction using only the trees for which that data point was not included in their respective bootstrap samples. This prediction is compared to the true label of the data point to compute the OOB error.
4. **Aggregating OOB Errors:** The OOB error for the entire random forest is computed by aggregating the errors over all data points in the original dataset. This aggregated error provides an estimate of how well the random forest will perform on unseen data.

The key advantage of using the out-of-bag error for model evaluation is that it provides an unbiased estimate of the model's performance without the need for a separate validation set. This makes the random forest algorithm computationally efficient and allows for more efficient use of the available data.

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

Answer K-fold cross-validation is a popular technique used to assess the performance and robustness of machine learning models. It involves partitioning the original dataset into k equal-

sized subsets, or folds. The model is then trained and evaluated  $k$  times, each time using a different fold as the validation set and the remaining  $k-1$  folds as the training set.

Here's how K-fold cross-validation works:

1. **Partitioning the Dataset:** The original dataset is divided into  $k$  equal-sized subsets, or folds. Each fold contains approximately  $n/k$  data points, where  $n$  is the total number of data points in the dataset.
2. **Training and Evaluation:** For each iteration  $i$  from 1 to  $k$ :
  - The model is trained on  $k-1$  folds (the training set) excluding the  $i$ th fold.
  - The trained model is then evaluated on the  $i$ th fold (the validation set) to calculate the performance metric(s) of interest, such as accuracy, precision, recall, or F1-score.
3. **Performance Metric Aggregation:** After performing K-fold cross-validation, the performance metrics obtained from each iteration are aggregated to provide a more robust estimate of the model's performance. This aggregation typically involves calculating the mean, median, or other summary statistics of the performance metrics across all  $k$  folds.

K-fold cross-validation is beneficial for several reasons:

- It provides a more reliable estimate of a model's performance compared to a single train-test split, as it uses multiple iterations and evaluates the model on different subsets of the data.
- It helps in detecting overfitting by assessing the model's performance on unseen data across multiple folds.
- It maximizes the use of the available data for both training and validation, which is especially useful for small datasets.

A common choice for the value of  $k$  is 5 or 10, but the optimal value depends on factors such as the size of the dataset and computational resources available.

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

Answer

Hyperparameter tuning in machine learning refers to the process of selecting the optimal hyperparameters for a given learning algorithm. Hyperparameters are parameters that are not directly learned from the training data but instead control the learning process. Examples of hyperparameters include the learning rate in gradient descent, the depth of a decision tree,

Hyperparameter tuning is done for several reasons:

1. **Optimizing Model Performance:** The choice of hyperparameters can significantly impact the performance of a machine learning model. By tuning hyperparameters, we aim to find the combination that results in the best possible performance on unseen data, such as higher accuracy, lower error, or improved generalization.
2. **Preventing Overfitting:** Hyperparameters control the complexity of a model, and improper settings can lead to overfitting, where the model learns to fit the training data too closely and fails to generalize well to new data. Tuning hyperparameters helps strike a balance between bias and variance, leading to a model that generalizes well to unseen data.
3. **Improving Interpretability:** In some cases, tuning hyperparameters can lead to models that are more interpretable and easier to understand. For example, reducing the maximum depth of a decision tree can lead to a simpler tree structure that is easier to interpret.
4. **Addressing Model Sensitivity:** Certain hyperparameters may make a model more sensitive to variations in the training data or the initial conditions of the optimization process. Tuning these hyperparameters can lead to models that are more robust and stable.

Hyperparameter tuning is typically done using techniques such as grid search, random search, Bayesian optimization, or more

advanced methods like genetic algorithms or reinforcement learning. These techniques systematically explore the hyperparameter space to find the combination that optimizes the chosen performance metric. Overall, hyperparameter tuning is a crucial step in the machine learning pipeline to ensure that models perform well and generalize effectively to new data.

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

Answer When using **Gradient Descent**, having a **large learning rate** can lead to several issues:

1. **Overshooting the Minimum:** A large learning rate can cause the updates to the model parameters to be too large, causing them to overshoot the minimum of the loss function. This can lead to the algorithm diverging, where the parameter values oscillate or diverge indefinitely, preventing convergence to the optimal solution..
2. **Instability:** When the learning rate is too large (e.g.,  $\text{learn\_rate} \geq 0.3$ ), the optimization process becomes unstable. Coefficients can explode, and overflow errors may occur. The algorithm takes excessively large steps, continuously missing the optimal solution.
3. **Instability and Unpredictable Behavior:** Large learning rates can introduce instability and erratic behavior in the optimization process. The algorithm may oscillate around the minimum or exhibit chaotic behavior, making it difficult to predict the convergence behavior.
4. **Divergence:** A large learning rate can cause the algorithm to diverge instead of converging to the minimum. Instead of gradually approaching the optimal point, it veers away, resulting in poor results.
5. **Slow or No Convergence:** In some cases, a large learning rate can prevent the algorithm from converging to the optimal solution altogether. Instead, it may bounce around the minimum or fail to make progress towards it, resulting in slow convergence or no convergence at all

In summary, choosing an appropriate learning rate is crucial for successful convergence in Gradient Descent. Too small a rate leads to slow convergence, while too large a rate can cause instability and divergence

To mitigate these issues, it is essential to choose an appropriate learning rate for gradient descent. This is often done through hyperparameter tuning or using techniques such as learning rate schedules, adaptive learning rate methods (e.g., Adam, RMSprop), or gradient clipping to prevent the gradients from becoming too large. Experimentation and monitoring of the optimization process are crucial to finding the optimal learning rate that balances convergence speed, stability, and generalization performance.

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

Answer

Logistic regression is a linear classification algorithm, meaning it assumes a linear relationship between the input features and the log-odds of the target variable. Therefore, logistic regression is not well-suited for handling non-linear relationships in the data.

When the relationship between the input features and the target variable is non-linear, logistic regression may not be able to capture the underlying patterns effectively. This can lead to poor performance and inaccurate predictions.

However, it's important to note that logistic regression can still be applied to non-linear data with some modifications, such as:

1. **Feature Engineering:** Transforming the input features or creating new features using non-linear transformations (e.g., polynomial features, interaction terms) can help capture non-linear relationships and make logistic regression more flexible.

2. **Kernel Methods:** Using kernel methods, such as kernel logistic regression or support vector machines with a non-linear kernel (e.g., polynomial kernel, radial basis function kernel), allows logistic regression to operate in a higher-dimensional feature space where non-linear relationships can be linearized.
3. **Ensemble Methods:** Combining multiple logistic regression models or using ensemble techniques like random forests or gradient boosting can also help handle non-linearities in the data by combining the strengths of multiple models.

In summary, while logistic regression is not inherently designed to handle non-linear data, it can still be used with appropriate modifications or in combination with other techniques to address non-linear relationships in the data. However, for highly non-linear data, more flexible models like neural networks or kernel methods may be more appropriate.

13. Differentiate between Adaboost and Gradient Boosting.

Answer

Adaboost and Gradient Boosting are both popular ensemble learning techniques used for building powerful predictive models, but they differ in their approach to training and updating the weak learners (base models) in the ensemble.

Here are the main differences between Adaboost and Gradient Boosting:

#### 1. Approach:

- **Adaboost (Adaptive Boosting):** Adaboost is an iterative ensemble technique that sequentially trains a sequence of weak learners (e.g., decision trees) on repeatedly modified versions of the data. It focuses on improving the performance of the weak learners by assigning higher weights to misclassified data points in each iteration.
- **Gradient Boosting:** Gradient Boosting is also an iterative ensemble technique that sequentially trains a sequence of weak learners (usually decision trees) on the residuals (the differences between the actual and predicted values) of the

previous learners. It focuses on optimizing the loss function by fitting subsequent models to the residuals of the previous models.

## 2. Weighting of Data Points:

- **Adaboost:** Adaboost assigns higher weights to misclassified data points in each iteration, thereby focusing more on the difficult-to-classify instances.
- **Gradient Boosting:** Gradient Boosting focuses on minimizing the errors (residuals) made by the previous models by fitting subsequent models to the residuals. It pays more attention to the data points with larger residuals.

## 3. Model Training:

- **Adaboost:** In each iteration, Adaboost selects the weak learner that best fits the data by minimizing the weighted error.
- **Gradient Boosting:** In each iteration, Gradient Boosting fits the weak learner to the residuals (errors) of the previous models. It updates the model parameters (e.g., tree structure, leaf values) in the direction that minimizes the loss function.

## 4. Loss Function:

- **Adaboost:** Adaboost uses exponential loss function for classification problems and exponential loss for regression problems.
- **Gradient Boosting:** Gradient Boosting allows flexibility in the choice of loss function and can be used with various loss functions such as squared loss (for regression), logistic loss (for binary classification), or softmax loss (for multi-class classification).

In summary, while both Adaboost and Gradient Boosting are ensemble techniques that combine multiple weak learners to build a strong predictive model, they differ in their approach to training and updating the weak learners. Adaboost focuses on sequentially adjusting the weights of the data points to improve model performance, while Gradient Boosting focuses on



sequentially fitting models to the residuals of the previous models to minimize the loss function.

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

Answer

The **bias-variance trade-off** is a fundamental concept in machine learning that affects the performance of predictive models. Let's break it down:

### **Bias:**

- **Bias** refers to the **difference between the predictions made by a machine learning model** and the actual (correct) values.
- When a model has **high bias**, it tends to make **simplistic assumptions** about the data, resulting in **underfitting**.
- Underfitting occurs when the model is too simple to capture the underlying patterns in the data.
- Imagine a straight-line hypothesis that doesn't fit the data well.
- **Goal:** Keep bias low to avoid underfitting

### **2. Variance:**

- **Variance** measures the **variability of model predictions** for a given data point.
- A model with **high variance** fits the training data very closely but performs poorly on unseen data (test data).
- High variance leads to **overfitting**, where the model captures noise and specific details of the training data.
- Overfitting results in poor generalization to new data.
- Think of a complex curve that perfectly fits the training points but fails to generalize.
- **Goal:** Keep variance low to avoid overfitting.

### **3. Trade-off:**

- The **bias-variance trade-off** arises because it's challenging to simultaneously minimize both bias and variance.
- **Complexity** plays a crucial role:

- If the model is too simple (low complexity), it has high bias and low variance.
- If the model is too complex (high complexity), it has low bias and high variance.
- The trade-off involves finding the right balance:
- **Optimal Model:** We seek a model that minimizes the **total error** (bias + variance) up to a certain point.
- This balance ensures good performance on both training and test data.
- **Graphical Representation:**
- The graph shows the trade-off region where the total error is minimized.
- Beyond that point, increasing complexity leads to overfitting.
- The sweet spot lies in the middle, where bias and variance are balanced.

Remember, achieving the right bias-variance trade-off is crucial for building models that generalize well to unseen data while avoiding both underfitting and overfitting

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

Answer

A short description of each of the kernels commonly used in Support Vector Machines (SVM):

### 1. **Linear Kernel:**

- The **linear kernel** is the simplest one. It creates a **linear decision boundary** in the feature space.
- It works well when the data is **linearly separable** or nearly so.
- The decision boundary is a **straight line** (or hyperplane in higher dimensions).
- It's computationally efficient and often used as a baseline.
- Example: Imagine drawing a straight line to separate two classes of data points.

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

- The **Radial Basis Function (RBF)** kernel is also known as the **Gaussian kernel**.
- It uses **normal curves** centered around data points to define the decision boundary.
- The boundary is determined by a **sum of these curves**, creating a flexible shape.
- It's effective for capturing intricate patterns and works well even when data is not linearly separable.
- Example: Picture a wavy boundary formed by combining many Gaussian distributions

## 3. Polynomial Kernel:

- The **polynomial kernel** introduces **non-linearity** by mapping the data into a higher-dimensional space.
- The boundary is defined by a **polynomial function** of some order (e.g., quadratic, cubic, etc.).
- It can capture more complex relationships between features.
- Example: Visualize a curved boundary, like a circle or an ellipse, separating classes.

Each of these kernels has its strengths and weaknesses, and the choice of kernel depends on the nature of the data and the problem at hand. Experimentation and cross-validation are often used to determine the most suitable kernel for a given task.