

# ASSIGNMENT - 5

## MACHINE LEARNING

**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 - R-squared is generally a better measure of the goodness of fit for a regression model than the residual sum of squares (RSS).

**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.**

Ans –

**Total Sum of Squares (TSS):**

- **Definition:** TSS measures the total variability of the observed data around the mean of the response variable yyy. It represents the total variation in the response variable that we are trying to explain.
- **Formula:**

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

**Explained Sum of Squares (ESS):**

- **Definition:** ESS measures the portion of the total variability that is explained by the regression model. It reflects how much of the variability in yyy is accounted for by the predictors.
- **Formula:**

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

**Residual Sum of Squares (RSS):**

**Definition:** RSS measures the portion of the total variability that is not explained by the regression model. It represents the error or residuals in the predictions.

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

**Formula:**

**Relationship between TSS, ESS, and RSS:**

The relationship between these three metrics is given by the following equation:

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

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

Ans –

#### 1. Prevent Overfitting

- **Problem:** Overfitting occurs when a model learns the noise and details in the training data to the extent that it negatively impacts its performance on unseen data. This typically happens when the model is too complex, with too many parameters relative to the amount of training data.
- **Solution:** Regularization techniques penalize the complexity of the model, discouraging it from fitting the training data too closely. This helps the model generalize better to new, unseen data.

#### 2. Simplify the Model

- **Problem:** Complex models with many parameters can become unwieldy and difficult to interpret. They may also be prone to learning irrelevant features or interactions that do not generalize well.
- **Solution:** Regularization often encourages simpler models by shrinking the coefficients of less important features towards zero. This can lead to a more interpretable and robust model.

#### 3. Improve Model Performance

- **Problem:** Without regularization, a model may achieve high performance on the training data but perform poorly on validation or test data due to overfitting.
- **Solution:** Regularization can improve performance on validation and test sets by reducing overfitting and ensuring that the model captures the underlying patterns rather than the noise.

#### 4. Handle Multicollinearity

- **Problem:** In regression models, multicollinearity occurs when predictors are highly correlated with each other, making it difficult to determine the individual effect of each predictor. This can lead to unstable estimates of coefficients.
- **Solution:** Regularization techniques like Ridge Regression (L2 regularization) can help mitigate the impact of multicollinearity by penalizing large coefficients and stabilizing the model.

## 5. Promote Sparse Solutions

- **Problem:** In high-dimensional settings, having many features can lead to overfitting and increased computational complexity.
- **Solution:** Regularization techniques like Lasso (L1 regularization) promote sparsity by shrinking some coefficients exactly to zero, effectively performing feature selection and reducing the dimensionality of the model.

## 4. What is Gini-impurity index?

Ans - The Gini impurity index, or Gini index, is a metric used in decision tree algorithms to measure the impurity or heterogeneity of a node. It helps in deciding how to split the data at each node of the tree to improve classification performance. Here's a detailed look at the Gini impurity index:

### Definition

The Gini impurity index is used to quantify the impurity of a dataset. It is commonly used in decision trees for classification tasks. The Gini impurity of a node is calculated as follows:

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

Ans - Yes, unregularized decision trees are indeed prone to overfitting. Here's why:

### Reasons for Overfitting in Unregularized Decision Trees

#### 1. High Model Complexity:

- **Issue:** Decision trees can grow very complex, especially if allowed to expand without constraints. A fully-grown tree might have many branches and nodes, each finely tuned to the training data.
- **Impact:** This high complexity enables the tree to capture not only the underlying patterns but also the noise and outliers in the training data. As a result, the model performs well on training data but poorly on unseen data.

#### 2. Over-Specialization:

- **Issue:** An unregularized decision tree splits nodes based on features in a way that reduces the impurity (e.g., Gini impurity or entropy) at each step. Without restrictions, the tree will continue splitting until each leaf node is as pure as possible, meaning it can end up overfitting to the training data.
- **Impact:** The result is a model that has memorized the training data rather than learning generalizable patterns, leading to poor generalization performance.

#### 3. Increased Variance:

- **Issue:** Decision trees have high variance when they are deep or unpruned. Small changes in the training data can lead to significant changes in the structure of the tree.
- **Impact:** This high variance means the model can be overly sensitive to the specifics of the training data, resulting in a lack of robustness when encountering new data.

#### 4. Lack of Generalization:

- **Issue:** Unregularized trees can make very specific decisions at each node, which are tailored to the exact values in the training set.
- **Impact:** As a consequence, these decisions might not generalize well to new or slightly different data, leading to overfitting.

### How Regularization Helps

To combat overfitting, several regularization techniques can be applied to decision trees:

#### 1. Pruning:

- **Description:** Pruning involves removing branches from the tree that have little importance, based on criteria such as the minimal contribution to reducing impurity.
- **Benefit:** This helps in simplifying the model and improving generalization by removing nodes that capture noise.

#### 2. Limiting Tree Depth:

- **Description:** Restricting the maximum depth of the tree prevents it from growing too complex.
- **Benefit:** This helps in controlling the complexity of the model and thus reduces overfitting.

#### 3. Minimum Samples Split and Leaf Size:

- **Description:** Setting minimum thresholds for the number of samples required to split a node or to form a leaf can prevent the tree from creating nodes with very few samples.
- **Benefit:** This prevents the tree from making decisions based on very small sample sizes, which can lead to overfitting.

#### 4. Cost Complexity Pruning (CCP):

- **Description:** Also known as weakest link pruning, this technique involves pruning the tree based on a complexity parameter that balances the tree's fit to the training data and its complexity.
- **Benefit:** It ensures that the trade-off between model fit and complexity is optimized.

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

Ans - Ensemble techniques in machine learning refer to methods that combine multiple models to improve the overall performance of a predictive system. The idea is that by aggregating the

predictions of multiple models, the ensemble can achieve better performance than any individual model alone. This is often because different models can capture different aspects of the data or make different types of errors, and combining them helps to average out these errors and reduce variance.

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

Ans –

### Bagging

#### 1. Purpose:

- **Goal:** Reduce variance and prevent overfitting by combining the predictions of multiple models trained on different subsets of the data.

#### 2. Training Process:

- **Data Sampling:** Creates multiple subsets of the training data by sampling with replacement (bootstrapping). Each model is trained on a different subset.
- **Model Training:** Models are trained independently on their respective subsets. Each model sees a different version of the training data.

#### 3. Aggregation:

- **Combining Predictions:** Predictions from the individual models are aggregated, typically using voting (for classification) or averaging (for regression).
- **Example:** Random Forest is a popular bagging technique where multiple decision trees are trained on different subsets of data and their predictions are averaged.

#### 4. Model Independence:

- **Independence:** Models are trained independently of each other, and there is no interaction between them during the training process.

#### 5. Effect on Bias and Variance:

- **Variance Reduction:** Bagging helps in reducing the variance of the model. It generally does not significantly affect the bias of the base models.

#### 6. Implementation:

- **Complexity:** Typically simpler to implement and understand compared to boosting.

### Boosting

#### 1. Purpose:

- **Goal:** Improve model performance by sequentially training models that focus on correcting the errors made by previous models. It reduces both bias and variance.

## 2. Training Process:

- **Sequential Learning:** Models are trained sequentially, with each new model attempting to correct the errors made by the previous models.
- **Weight Adjustment:** The training data's weights are adjusted based on the errors made by previous models. Misclassified examples are given more weight to focus on correcting them.

## 3. Aggregation:

- **Combining Predictions:** Predictions are typically combined using weighted voting or averaging. Each model's contribution is weighted based on its performance.
- **Example:** AdaBoost and Gradient Boosting Machines (GBM) are popular boosting techniques. In AdaBoost, for instance, each subsequent model focuses on the mistakes made by the previous models.

## 4. Model Dependency:

- **Dependency:** Models are not trained independently; each model is dependent on the previous one. The sequence of training builds on the mistakes of prior models.

## 5. Effect on Bias and Variance:

- **Bias and Variance Reduction:** Boosting can reduce both bias and variance. It tends to reduce bias by focusing on difficult examples, which often leads to improved overall accuracy.

## 6. Implementation:

- **Complexity:** Typically more complex to implement and tune due to the sequential nature and adjustments needed for weights.

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

Ans - The out-of-bag error in Random Forests is a practical and efficient way to estimate the model's prediction error. It uses the instances not included in the bootstrap sample for each tree to test the model, providing a built-in cross-validation process that helps gauge the model's performance without requiring a separate validation set. This method contributes to the robustness of Random Forests by providing a reliable estimate of how well the ensemble of trees is likely to perform on unseen data.

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

Ans - K-fold cross-validation is a robust method for evaluating a machine learning model's performance. By systematically partitioning the dataset into K folds and using each fold for validation while training on the remaining folds, it provides a comprehensive estimate of how the model will perform on unseen data. This approach helps in making the best use of available data and provides a more reliable assessment of model performance.

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

Ans **Hyperparameter tuning** in machine learning refers to the process of selecting the best set of hyperparameters for a given model. Hyperparameters are parameters that are set before the learning process begins and are not learned from the data during training. They control various aspects of the model's training process and architecture, influencing how well the model performs.

### What Are Hyperparameters?

- **Definition:** Hyperparameters are settings or configurations external to the model that influence the learning process. They are different from model parameters, which are learned from the data.
- **Examples:**
  - **Learning Rate:** In gradient-based algorithms, it controls the size of the steps taken towards the minimum of the loss function.
  - **Number of Trees:** In Random Forests, it specifies how many decision trees are built.
  - **Depth of Trees:** In decision trees and Random Forests, it specifies how deep the trees should be.
  - **Number of Hidden Layers and Units:** In neural networks, these specify the architecture of the network.
  - **Regularization Parameters:** Such as  $\lambda$  in Ridge regression or  $\alpha$  in Lasso regression.

### Why Hyperparameter Tuning Is Done

#### 1. Improves Model Performance:

- **Optimal Settings:** Choosing the best hyperparameters can significantly enhance the model's performance, including accuracy, precision, recall, and other evaluation metrics.
- **Model Effectiveness:** Properly tuned hyperparameters help in better capturing the underlying patterns in the data and reduce errors.

#### 2. Prevents Overfitting and Underfitting:

- **Balancing Act:** Hyperparameter tuning helps in finding the right balance between underfitting (model is too simple) and overfitting (model is too complex). For instance, tuning the complexity of the model can help in generalizing better to unseen data.

#### 3. Optimizes Computational Resources:

- **Efficiency:** Proper tuning can lead to a more efficient use of computational resources by avoiding excessively complex models that may be computationally expensive.

#### 4. Adapts to Different Data:

- **Flexibility:** Different datasets may require different hyperparameter settings for optimal performance. Tuning ensures that the model is well-suited for the specific characteristics of the data.

### Methods for Hyperparameter Tuning

#### 1. Grid Search:

- **Description:** A systematic approach where a predefined set of hyperparameter values is specified, and the model is trained and evaluated for every combination of these values.
- **Advantage:** Exhaustive and ensures that all possible combinations within the grid are evaluated.
- **Disadvantage:** Computationally expensive, especially with a large number of hyperparameters and their possible values.

#### 2. Random Search:

- **Description:** Randomly samples a subset of hyperparameter combinations from a predefined range.
- **Advantage:** More efficient than grid search as it explores a wider range of values with less computational effort.
- **Disadvantage:** May miss the optimal combination if it is not sampled.

#### 3. Bayesian Optimization:

- **Description:** A probabilistic model-based approach that builds a surrogate model to predict the performance of different hyperparameter combinations and uses it to select the next set of hyperparameters to evaluate.
- **Advantage:** More efficient in exploring the hyperparameter space by focusing on promising areas.
- **Disadvantage:** Can be complex to implement and computationally intensive.

#### 4. Automated Machine Learning (AutoML):

- **Description:** Tools and frameworks that automate the process of hyperparameter tuning and model selection.
- **Advantage:** Simplifies the process and can be very effective in finding good hyperparameters.
- **Disadvantage:** May not provide the same level of control or insight into the tuning process.



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

Ans - In gradient descent, the learning rate is a crucial hyperparameter that controls how much to adjust the model parameters in response to the estimated gradient of the loss function. If the learning rate is too large, several issues can arise:

### Issues with a Large Learning Rate

#### 1. Overshooting the Minimum:

- **Description:** A large learning rate can cause the algorithm to take steps that are too big, causing it to overshoot the minimum of the loss function. Instead of converging smoothly towards the minimum, the parameter updates might jump over it.
- **Effect:** This can lead to oscillations around the minimum or even divergence, where the model parameters continue to move further away from the optimal point.

#### 2. Divergence:

- **Description:** In extreme cases, a learning rate that is too large can lead to divergence, where the model's loss increases instead of decreasing. This happens because the parameter updates are so large that they move the model parameters far from the optimal values, causing the loss to grow.
- **Effect:** Divergence means the algorithm fails to converge and may produce an unstable model that cannot be used for predictions.

#### 3. Instability:

- **Description:** A large learning rate can cause the optimization process to become unstable. The loss function may fluctuate wildly rather than steadily decreasing, making it difficult to reach a satisfactory solution.
- **Effect:** This instability makes it challenging to monitor and control the training process, potentially leading to a model that does not generalize well.

#### 4. Poor Convergence:

- **Description:** Even if the gradient descent does not completely diverge, a large learning rate can still lead to poor convergence. The algorithm may take many iterations to settle down to the optimal values or get stuck in local minima.
- **Effect:** This inefficiency can result in longer training times and suboptimal model performance.

### Visualizing the Problem

#### 1. Loss Function Behavior:

- **Plot:** When the learning rate is too large, the plot of the loss function over iterations often shows large jumps or oscillations rather than a smooth decline.

#### 2. Parameter Trajectories:

- **Plot:** The trajectory of the model parameters in parameter space can appear erratic or zigzag, failing to settle into a stable region near the minimum.

### Mitigating the Issues

#### 1. Learning Rate Scheduling:

- **Description:** Use techniques to reduce the learning rate over time, such as step decay, exponential decay, or adaptive learning rates (e.g., using optimizers like Adam or RMSprop).
- **Benefit:** This helps in taking larger steps initially when far from the minimum and smaller steps as the algorithm converges closer to the optimal point.

#### 2. Learning Rate Annealing:

- **Description:** Start with a larger learning rate and gradually decrease it during training.
- **Benefit:** Allows the algorithm to explore the parameter space more broadly at first and then fine-tune the parameters with a smaller learning rate.

#### 3. Using Optimizers with Adaptive Learning Rates:

- **Description:** Algorithms like Adam, RMSprop, or AdaGrad adjust the learning rate dynamically based on the gradient history or other criteria.
- **Benefit:** These optimizers can adaptively control the learning rate to prevent issues with too large updates.

#### 4. Grid Search or Random Search for Learning Rate:

- **Description:** Experiment with different learning rates using grid search or random search to find an optimal value.
- **Benefit:** Helps in finding a learning rate that balances convergence speed and stability.

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

Ans - Logistic Regression is primarily designed for binary classification tasks and works best for linearly separable data. When dealing with non-linear data, where the decision boundary between classes is not a straight line, Logistic Regression by itself may not be suitable. Here's why and what can be done to address this limitation:

### Why Logistic Regression Struggles with Non-Linear Data

#### 1. Linear Decision Boundary:

- **Fundamental Issue:** Logistic Regression assumes that the relationship between the features and the probability of the target class is linear. It models the probability of a binary outcome as a logistic function (sigmoid function) of a linear combination of input features.

- **Implication:** This results in a linear decision boundary in the feature space, which means it can only separate classes with a straight line (or hyperplane in higher dimensions).

## 2. Limited Flexibility:

- **Model Limitation:** Because of its linear nature, Logistic Regression is not inherently capable of capturing more complex, non-linear relationships between the features and the target variable.

## How to Handle Non-Linear Data with Logistic Regression

### 1. Feature Engineering:

- **Polynomial Features:** Transform the original features into polynomial features (e.g.,  $x_1^2$ ,  $x_1 \cdot x_2$ ) to allow the model to fit more complex, non-linear relationships.
- **Interaction Terms:** Include interaction terms between features to capture non-linear interactions.

### 2. Kernel Trick:

- **Extension:** Use a kernel function to implicitly map the input features into a higher-dimensional space where a linear decision boundary might be more appropriate. This approach is more commonly used with Support Vector Machines (SVMs) rather than Logistic Regression but illustrates the principle of using a higher-dimensional space for non-linear separation.

### 3. Non-Linear Transformations:

- **Feature Transformations:** Apply non-linear transformations to the features (e.g., log, square root) to make the relationship more linear in the transformed feature space.

### 4. Extended Models:

- **Polynomial Logistic Regression:** Involves adding polynomial terms of the features to the model to capture non-linearity. This is effectively a form of feature engineering that makes the logistic model more flexible.
- **Non-Linear Models:** Use models specifically designed for non-linear classification, such as:
  - **Decision Trees:** Can model non-linear relationships with their hierarchical structure.
  - **Random Forests:** Ensembles of decision trees that can capture complex non-linear patterns.
  - **Support Vector Machines (SVMs):** With kernel functions, SVMs can handle non-linear data.
  - **Neural Networks:** Highly flexible and capable of learning complex non-linear relationships.

## 13. Differentiate between Adaboost and Gradient Boosting.

Ans - AdaBoost (Adaptive Boosting) and Gradient Boosting are both boosting techniques used to improve the performance of machine learning models by combining multiple weak learners (usually decision trees) into a single strong model. However, they have different approaches and mechanisms. Here's a detailed differentiation between AdaBoost and Gradient Boosting:

### AdaBoost

#### 1. Basic Principle:

- **Focus on Misclassified Examples:** AdaBoost works by sequentially training weak learners, where each new learner focuses on correcting the errors made by the previous ones. It adjusts the weights of incorrectly classified examples to emphasize them more in subsequent models.

#### 2. Training Process:

- **Sequential Adjustment:** Each new model is trained on a weighted version of the data. Initially, all examples have equal weights, but after each model is trained, the weights of misclassified examples are increased while correctly classified examples have their weights decreased.
- **Final Model:** The predictions from all models are combined through weighted voting, where each model's vote is weighted by its accuracy.

#### 3. Algorithm Steps:

1. Initialize weights for all training examples.
2. Train a weak learner on the weighted data.
3. Compute the error of the weak learner.
4. Update the weights of the examples based on the errors made.
5. Combine the predictions of all weak learners, with more weight given to those with lower error rates.

#### 4. Loss Function:

- **Exponentiated Loss Function:** AdaBoost minimizes an exponential loss function, which emphasizes incorrect predictions more.

#### 5. Robustness:

- **Sensitivity to Noise:** AdaBoost can be sensitive to noisy data and outliers because it focuses heavily on correcting misclassified examples.

#### 6. Implementation:

- **Example:** AdaBoost can be implemented with a variety of base learners, though decision stumps (one-level decision trees) are common.

### Gradient Boosting

#### 1. Basic Principle:

- **Gradient Descent on Loss Function:** Gradient Boosting builds models sequentially where each new model tries to correct the errors of the combined ensemble of previous models by minimizing a specific loss function through gradient descent.

## 2. Training Process:

- **Additive Model:** Each new model is trained to predict the residual errors (the difference between the actual values and the predictions of the ensemble of previously trained models).
- **Final Model:** The predictions from all models are combined by adding the outputs of each weak learner.

## 3. Algorithm Steps:

1. Initialize the model with a simple base model (often a constant prediction).
2. Compute the residuals (errors) between the true values and the model's predictions.
3. Train a new weak learner to predict these residuals.
4. Update the model by adding the predictions of the new weak learner.
5. Repeat the process for a set number of iterations or until no further improvement is observed.

## 4. Loss Function:

- **Flexible Loss Functions:** Gradient Boosting can use various loss functions (e.g., mean squared error for regression, log loss for classification) and is based on gradient descent optimization.

## 5. Robustness:

- **Flexibility:** Gradient Boosting is more flexible and can handle different types of loss functions, making it suitable for a variety of problems. It is generally less sensitive to noise than AdaBoost.

## 6. Implementation:

- **Example:** Gradient Boosting implementations include popular libraries such as XGBoost, LightGBM, and CatBoost, which offer optimizations and enhancements over basic gradient boosting.

## Summary of Differences

### 1. Approach:

- **AdaBoost:** Focuses on correcting errors of previous models by adjusting weights of misclassified examples.
- **Gradient Boosting:** Optimizes a loss function through gradient descent by fitting new models to the residuals of the combined ensemble.

### 2. Loss Function:

- **AdaBoost:** Uses an exponential loss function.

- **Gradient Boosting:** Can use various loss functions based on the specific problem (e.g., mean squared error, log loss).

### 3. **Model Adjustment:**

- **AdaBoost:** Adjusts weights of training examples to focus on harder-to-classify examples.
- **Gradient Boosting:** Adds new models to minimize residuals of the existing ensemble.

### 4. **Sensitivity to Noise:**

- **AdaBoost:** Can be more sensitive to noisy data and outliers.
- **Gradient Boosting:** More flexible and generally less sensitive to noise.

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

Ans - The bias-variance trade-off is a crucial aspect of model development in machine learning. It involves finding the right balance between bias (error due to overly simplistic assumptions) and variance (error due to model complexity and sensitivity to data fluctuations). Managing this trade-off effectively helps in building models that generalize well to new, unseen data, avoiding both underfitting and overfitting.

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

Ans –

**Linear Kernel:** Suitable for linearly separable data; computationally simple.

**Polynomial Kernel:** Handles non-linear data by using polynomial features; flexibility depends on the polynomial degree.

**RBF Kernel:** Handles complex, non-linear data with a flexible decision boundary; highly effective for various data distributions.