

## STATISTICS WORKSHEET - 5

1. Using a goodness of fit, we can assess whether a set of obtained frequencies differ from a set of frequencies. a) Mean b) Actual c) Predicted d) Expected

Answer: d)

2. Chi-square is used to analyse a) Score b) Rank c) Frequencies d) All of these

Answer: c)

3. What is the mean of a Chi Square distribution with 6 degrees of freedom? a) 4 b) 12 c) 6 d) 8

Answer: c)

4. Which of these distributions is used for a goodness of fit testing? a) Normal distribution b) Chi-squared distribution c) Gamma distribution d) Poisson distribution

Answer: b)

5. Which of the following distributions is Continuous a) Binomial Distribution b) Hypergeometric Distribution c) F Distribution d) Poisson Distribution

Answer: c)

6. A statement made about a population for testing purpose is called? a) Statistic b) Hypothesis c) Level of Significance d) Test Statistic

Answer: b)

7. If the assumed hypothesis is tested for rejection considering it to be true is called? a) Null Hypothesis b) Statistical Hypothesis c) Simple Hypothesis d) Composite Hypothesis

Answer: a)

8. If the Critical region is evenly distributed then the test is referred as? a) Two tailed b) One tailed c) Three tailed d) Zero tailed

Answer: a)

9. Alternative Hypothesis is also called as? a) Composite hypothesis b) Research Hypothesis c) Simple Hypothesis d) Null Hypothesis

Answer b)

10. In a Binomial Distribution, if 'n' is the number of trials and 'p' is the probability of success, then the mean value is given by a) np b) n

Answer a)

## MACHINE LEARNING ASSIGNMENT - 5

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:

In regression analysis, both R-squared ( $R^2$ ) and Residual Sum of Squares (RSS) are used to measure the goodness of fit of a model, but they convey this information in different ways. Understanding which one is better depends on the context and the specific aspects of model performance you are focusing on.

### Comparison and Context

- **Ease of Interpretation:** R-squared is generally easier to interpret because it provides a standardized measure (0 to 1) of how well the model explains the variance. It allows for easier comparison between different models and is intuitively understood as a percentage of explained variance.
- **Absolute Measure:** RSS provides an absolute measure of fit, which can be useful in contexts where the actual magnitude of the errors is important. However, RSS is scale-dependent, meaning it is harder to compare between different datasets or models with different dependent variables.

- **Comparing Models:** If you are comparing multiple models with the same dependent variable, R-squared is often more informative as it standardizes the goodness of fit. However, when comparing models on different scales, RSS may not be as useful without additional context.

## Conclusion

**R-squared ( $R^2$ )** is generally a better measure for goodness of fit in regression when comparing models and understanding the proportion of explained variance. It is standardized, making it easier to interpret and compare across different models. However, **Residual Sum of Squares (RSS)** can be useful in specific contexts where the actual magnitude of residuals is important, but it is less intuitive and more challenging to use for comparison purposes across different datasets or models.

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, the Total Sum of Squares (TSS), Explained Sum of Squares (ESS), and Residual Sum of Squares (RSS) are fundamental metrics used to measure different aspects of the variability in the data. Here's a detailed explanation of each, along with the equation relating them.

### Total Sum of Squares (TSS)

- **Definition:** TSS measures the total variability in the dependent variable. It represents the total sum of the squared differences between the observed values and the mean of the observed values.
- **Formula:**  $TSS = \sum (y_i - \bar{y})^2$  where  $y_i$  is the observed values,  $\bar{y}$  is the mean of the observed values

### Explained Sum of Squares (ESS)

- **Definition:** ESS measures the variability explained by the regression model. It represents the sum of the squared differences between the predicted values and the mean of the observed values.
- **Formula:**  $ESS = \sum (y_x - \bar{y})^2$  where  $y_x$  are the predicted values of  $y$  from the regression model

### Residual Sum of Squares (RSS)

- **Definition:** RSS measures the variability that is not explained by the regression model. It represents the sum of the squared differences between the observed values and the predicted values.
- **Formula:**  $RSS = \sum (y_i - y_x)^2$  where  $y_i$  is the observed values,  $y_x$  are the predicted values

### Relationship Between TSS, ESS, and RSS

The Total Sum of Squares (TSS) is the sum of the Explained Sum of Squares (ESS) and the Residual Sum of Squares (RSS). This relationship can be expressed as:

$$TSS = ESS + RSS$$

This equation highlights that the total variability in the data (TSS) is partitioned into the variability explained by the model (ESS) and the unexplained variability (RSS).

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

**Answer:** Regularization in machine learning is a crucial technique used to prevent overfitting and improve the generalization of models. Overfitting occurs when a model learns not only the underlying pattern in the training data but also the noise, making it perform well on training data but poorly on unseen test data. Regularization helps to address this issue by adding a penalty term to the loss function, which discourages the model from becoming too complex.

#### Key Benefits and Needs for Regularization:

1. **Prevents Overfitting:**
  - Regularization adds a complexity penalty to the loss function, which discourages the model from fitting the training data too closely. This leads to a simpler model that is less likely to overfit the data.
2. **Improves Generalization:**
  - By penalizing overly complex models, regularization helps ensure that the model performs well on new, unseen data, improving its ability to generalize.
3. **Controls Model Complexity:**
  - Regularization techniques introduce additional terms that control the complexity of the model parameters, ensuring that the model remains robust and interpretable.

#### Summary:

Regularization is essential in machine learning for the following reasons:

- **Prevents Overfitting:** By penalizing large coefficients, it discourages the model from fitting the noise in the training data.
- **Improves Generalization:** Leads to better performance on unseen data by keeping the model simpler.
- **Controls Complexity:** Ensures that the model is not overly complex and is more interpretable.

By incorporating regularization techniques like L1, L2, or Elastic Net, machine learning models can achieve a good balance between bias and variance, leading to more reliable and robust predictions.

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

**Answer:** The Gini impurity index is a metric used in decision tree algorithms to measure the impurity or heterogeneity of a dataset. It indicates how often a randomly chosen element from the dataset would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the dataset. The Gini impurity is used to decide the optimal split at each node in the tree.

#### Definition

The Gini impurity for a dataset is defined as:

$Gini(D) = 1 - \sum p_j^2$  where  $p_j$  is the proportion of instances in the dataset that belong to class  $j$

### Interpretation

- **Gini Impurity of 0:** Indicates that all elements belong to a single class, meaning the dataset is perfectly pure.
- **Higher Gini Impurity:** Indicates a higher degree of impurity, with a more even distribution of elements among different classes.

### Use in Decision Trees

In decision tree algorithms, such as CART (Classification and Regression Trees), the Gini impurity is used to evaluate splits at each node. The goal is to select splits that reduce the impurity, leading to nodes that are as pure as possible. The split that results in the largest decrease in Gini impurity is chosen.

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

### Answer:

Yes, unregularized decision trees are indeed prone to overfitting.

### High Variance and Overfitting

1. **Flexibility of Decision Trees:**
  - Decision trees are highly flexible and can create very complex models by splitting the data into increasingly smaller subsets based on the features. This flexibility allows them to perfectly fit the training data, capturing all patterns, including noise and outliers.
2. **Depth of the Tree:**
  - Without regularization, decision trees can grow to an arbitrary depth, with many nodes and leaves. Each leaf can potentially contain only a single observation, leading to a model that fits the training data extremely well but fails to generalize to unseen data. This deep tree structure makes the model highly sensitive to the specific training data, resulting in high variance.
3. **Pure Nodes:**
  - An unregularized decision tree will continue splitting until it achieves pure nodes (where all samples in a node belong to a single class). While this might reduce impurity within the training set, it usually reflects overfitting, as the model captures noise specific to the training set.

### Lack of Regularization

1. **No Pruning:**
  - Regularization techniques such as pruning (pre-pruning or post-pruning) are not applied. Pre-pruning stops the tree from growing once it reaches a certain depth or number of leaves, while post-pruning removes branches that have little importance

after the tree is fully grown. Without pruning, the tree is likely to become overly complex.

2. **No Constraints:**

- Unregularized trees lack constraints like a minimum number of samples required to split a node, a maximum depth, or a minimum impurity decrease required for splitting. These constraints help in simplifying the tree and reducing the risk of overfitting.

6. What is an ensemble technique in machine learning?

**Answer:**

Ensemble techniques in machine learning involve combining multiple models to improve the overall performance, robustness, and generalization of the predictive model. By leveraging the strengths of multiple models, ensemble methods can achieve better results than any individual model alone. These techniques are particularly effective in reducing variance, bias, and improving prediction accuracy.

**Key Ensemble Techniques**

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

- **Description:** Bagging involves training multiple instances of the same model on different subsets of the training data. These subsets are generated by bootstrapping, where each subset is created by sampling with replacement.
- **Example:** Random Forests are a popular bagging technique where multiple decision trees are trained on bootstrapped subsets of the data, and their predictions are aggregated (usually by averaging for regression or voting for classification).

2. **Boosting:**

- **Description:** Boosting involves training multiple models sequentially, where each new model focuses on correcting the errors made by the previous models. The models are trained in a way that each subsequent model gives more weight to the instances that were previously misclassified.
- **Example:** Gradient Boosting, AdaBoost, and XGBoost are popular boosting techniques.

3. **Stacking (Stacked Generalization):**

- **Description:** Stacking involves training multiple base models (level-0 models) and then using their predictions as input features to a higher-level model (level-1 model) that makes the final prediction. The higher-level model learns how to best combine the base models' predictions.
- **Example:** Using different classifiers as base models and a logistic regression as the meta-model.
- **Description:** Voting involves training multiple models and combining their predictions by voting. In a classification problem, the final prediction is based on the majority vote (hard voting) or the average of the predicted probabilities (soft voting).
- **Example:** Combining the predictions of different classifiers like logistic

## Benefits of Ensemble Techniques

1. **Improved Accuracy:** By combining the strengths of multiple models, ensembles typically outperform individual models in terms of prediction accuracy.
2. **Robustness:** Ensembles are more robust to overfitting and variance, as the errors of individual models tend to cancel each other out.
3. **Stability:** Ensembles provide more stable and reliable predictions, reducing the impact of any single model's weaknesses.

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

**Answer:** Bagging (Bootstrap Aggregating) and Boosting are both ensemble techniques used in machine learning to improve the accuracy and robustness of predictive models. However, they have different approaches and characteristics. Here's a detailed comparison:

### Bagging (Bootstrap Aggregating)

1. **Purpose:**
  - Reduces variance and helps prevent overfitting.
2. **Method:**
  - Multiple models (usually of the same type) are trained in parallel on different subsets of the training data, generated by bootstrapping (sampling with replacement).
  - The final prediction is made by aggregating the predictions of all models (e.g., by averaging for regression or majority voting for classification).
3. **Characteristics:**
  - **Parallel Training:** All models are trained independently and in parallel.
  - **Model Independence:** Each model is built independently of the others.
  - **Variance Reduction:** By averaging multiple predictions, it reduces the model's variance and improves generalization.
4. **Example Algorithm:**
  - **Random Forest:** An ensemble of decision trees where each tree is trained on a bootstrapped sample of the data, and the final prediction is the average or majority vote of all trees.

### Boosting

1. **Purpose:**
  - Reduces both bias and variance, and helps improve model accuracy.
2. **Method:**
  - Multiple models are trained sequentially, where each new model focuses on correcting the errors made by the previous models.
  - The models are not built independently; instead, each subsequent model is influenced by the performance of the previous ones.
3. **Characteristics:**
  - **Sequential Training:** Models are trained one after another, with each new model attempting to correct the errors of the previous models.
  - **Weight Adjustment:** In many boosting algorithms, weights are adjusted to give more importance to the misclassified samples.

- **Bias and Variance Reduction:** By focusing on difficult cases and correcting errors, it reduces both bias and variance.
4. **Example Algorithms:**
- AdaBoost: Adjusts the weights of incorrectly classified samples so that subsequent models focus more on them.
  - Gradient Boosting: Builds models sequentially, with each new model trying to reduce the residual errors of the combined previous models.
  - XGBoost: An optimized and regularized version of Gradient Boosting that performs well in practice and is widely used in competitions.

**Comparison Table**

Feature	Bagging	Boosting
<b>Training Method</b>	Parallel (independent models)	Sequential (dependent models)
<b>Focus</b>	Reducing variance	Reducing bias and variance
<b>Model Interaction</b>	Independent models	Each model corrects previous errors
<b>Examples</b>	Random Forest	AdaBoost, Gradient Boosting, XGBoost
<b>Final Prediction</b>	Averaging (regression) or Voting (classification)	Weighted sum of predictions
<b>Data Subsets</b>	Bootstrapped subsets (with replacement)	Entire dataset, weighted samples

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

**Answer:** Out-of-bag (OOB) error is an internal method of measuring the prediction error of random forests without the need for a separate validation set. It leverages the bootstrap sampling technique used in constructing the random forest.

1. **Bootstrap Sampling:** When building each tree in a random forest, a bootstrap sample is created by randomly sampling the training data with replacement. Typically, about 63% of the original training data is used to train each tree, and the remaining 37% is not included in the sample.
2. **Out-of-Bag Data:** The 37% of the data that is not included in the bootstrap sample for a particular tree is referred to as the "out-of-bag" data for that tree.
3. **Prediction and Error Calculation:**
  - **Prediction:** Each tree in the random forest can make predictions on its own OOB data.
  - **Error Calculation:** The OOB error is calculated by aggregating the predictions for each data point that was not included in the bootstrap sample. Specifically, for each data point, you average the predictions from all trees that did not use that point in



their bootstrap sample. The error rate is then computed based on these aggregated predictions.

#### Advantages of OOB Error:

1. **No Need for a Separate Validation Set:** The OOB error provides an unbiased estimate of the model's performance without requiring a separate validation set. This is particularly useful when the dataset is small, as it allows all the data to be used for training while still providing a reliable measure of model performance.
2. **Unbiased Estimate:** OOB error is generally an unbiased estimate of the true prediction error, as it effectively uses a form of cross-validation within the random forest framework.

The out-of-bag (OOB) error is a useful method for estimating the prediction error of random forests. It takes advantage of the bootstrap sampling process inherent in random forests to provide an unbiased estimate of model performance without needing a separate validation set. This makes it a valuable tool for model evaluation, especially when data is limited.

9. What is K-fold cross-validation?

#### Answer:

K-fold cross-validation is a robust technique used in machine learning to evaluate the performance of a model and to ensure that it generalizes well to unseen data. It involves dividing the dataset into K equal-sized subsets, or "folds," and then performing the following steps:

#### Steps in K-fold Cross-validation

1. **Splitting the Data:**
  - The dataset is randomly partitioned into K equal-sized subsets (folds).
2. **Training and Validation:**
  - For each fold, the model is trained on  $K-1$  folds and validated on the remaining one fold.
  - This process is repeated K times, with each fold being used exactly once as the validation data.
3. **Averaging the Results:**
  - The performance metrics (e.g., accuracy, precision, recall, etc.) are computed for each of the K iterations.
  - The final performance metric is obtained by averaging the results from the K folds.

#### Advantages of K-fold Cross-validation

1. **More Reliable Estimates:**
  - By using multiple folds, K-fold cross-validation provides a more reliable and stable estimate of model performance compared to a single train-test split.
2. **Efficient Use of Data:**
  - Each data point is used both for training and validation, maximizing the utility of the dataset.
3. **Reduced Variance:**

- Since the model is evaluated on multiple folds, the variance in the performance metric is reduced, providing a better understanding of the model's performance on unseen data.

### Choosing the Value of K

- Common values of K are 5 and 10, but the choice can vary depending on the size of the dataset and the specific application.
- A larger K means that each training set will be more similar to the original dataset (less biased), but it also increases computational cost.

### Summary

K-fold cross-validation is a powerful technique for assessing the performance of machine learning models. It helps to mitigate overfitting, provides a more comprehensive evaluation of model performance, and ensures that the model generalizes well to unseen data. By training and validating the model on multiple subsets of the data, K-fold cross-validation provides more reliable and stable performance metrics.

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

**Answer:** Hyperparameter tuning in machine learning involves selecting the optimal set of hyperparameters for a learning algorithm. Hyperparameters are parameters that are not learned from the data but are set prior to the training process and govern the training behavior of the algorithm.

### Need:

1. **Model Performance:** Proper hyperparameter tuning can significantly improve the performance of a machine learning model. It helps in finding the best set of hyperparameters that minimize the error and maximize the predictive accuracy on unseen data.
2. **Generalization:** Hyperparameter tuning helps in building models that generalize well to new data, avoiding both underfitting and overfitting.
3. **Model Complexity:** Tuning can help in managing the complexity of the model by controlling parameters that determine the capacity of the learning algorithm, such as the depth of a decision tree or the number of hidden layers in a neural network.

### Common Hyperparameters

- **Learning Rate:** Controls how much the model's weights are adjusted with respect to the loss gradient.
- **Number of Trees (in Random Forests):** Determines the number of decision trees in an ensemble.
- **Depth of Tree (in Decision Trees and Random Forests):** Controls the maximum depth of the trees.
- **Number of Neighbors (in k-NN):** Specifies the number of nearest neighbors to consider.

- **Regularization Parameter:** Controls the amount of regularization applied to prevent overfitting (e.g., L2 penalty in Ridge Regression).
- **Batch Size:** Determines the number of samples used in one iteration to update the model parameters.

Hyperparameter tuning is a critical step in the machine learning pipeline that involves selecting the best hyperparameters to optimize model performance. By systematically searching through different combinations of hyperparameters, we can significantly improve the accuracy and generalization of the model. Various methods such as grid search, random search, and Bayesian optimization provide different approaches to efficiently explore the hyperparameter space.

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

**Answer:**

If the learning rate in Gradient Descent is set too large, several issues can occur that can hinder the training process and the performance of the machine learning model. Here are the main problems associated with a large learning rate:

1. **Overshooting the Minimum:**
  - **Description:** A large learning rate can cause the algorithm to overshoot the minimum of the cost function. Instead of gradually approaching the minimum, the updates to the model parameters are too large, causing the algorithm to jump back and forth across the minimum, never converging.
  - **Impact:** This results in failure to find the optimal solution and the training process oscillates or diverges.
2. **Divergence:**
  - **Description:** If the learning rate is excessively large, the updates can become so substantial that the model parameters move further away from the optimal point with each iteration, leading to an increase in the cost function.
  - **Impact:** The cost function values increase instead of decrease, causing the algorithm to diverge, and the training process becomes unstable.
3. **Poor Convergence:**
  - **Description:** Even if the learning rate is not large enough to cause complete divergence, it can still hinder convergence. The algorithm might take a very erratic path towards the minimum, causing slow and inefficient convergence.
  - **Impact:** This can result in a much longer training time and potentially never reaching the minimum within a reasonable number of iterations.
4. **High Variance in Cost Function:**
  - **Description:** With a large learning rate, the updates to the cost function can vary significantly, resulting in high variance. The cost function plot will show a lot of fluctuation, making it difficult to diagnose the training process.
  - **Impact:** This makes it challenging to monitor the progress of the training and to determine whether the algorithm is converging.

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

**Answer:** Logistic regression is inherently a linear model, which means it works well for data that is linearly separable. However, it struggles with non-linear data because it tries to separate the classes with a linear decision boundary. Here's why logistic regression alone is not suitable for non-linear data and some methods to address this limitation:

### Logistic Regression Struggles with Non-Linear Data

1. **Linear Decision Boundary:**
  - Logistic regression models the probability of a class label as a linear function of the input features. This results in a linear decision boundary in the feature space.
  - For non-linearly separable data, a linear decision boundary cannot effectively separate the classes, leading to poor classification performance.
2. **Feature Space Limitation:**
  - Logistic regression directly uses the input features as they are. If the relationship between the features and the target variable is non-linear, the model cannot capture this complexity.

### Addressing Non-Linearity

While logistic regression itself is linear, there are several ways to make it handle non-linear data:

1. **Polynomial Features:**
  - Transform the input features into polynomial features to introduce non-linearity. This approach allows logistic regression to fit a non-linear decision boundary by using the original features in a non-linear manner.
2. **Interaction Terms:**
  - Include interaction terms between features to capture their combined effect. Interaction terms can help capture complex relationships in the data.
3. **Kernel Methods:**
  - Use kernel tricks, similar to support vector machines (SVMs), to map the original features into a higher-dimensional space where a linear decision boundary can be effective. However, logistic regression does not natively support kernel methods, so this approach often involves using SVMs instead.
4. **Non-Linear Models:**
  - Consider using inherently non-linear models such as decision trees, random forests, gradient boosting machines, neural networks, or support vector machines with non-linear kernels (e.g., RBF kernel).

While logistic regression is a linear model and struggles with non-linear data, it can be extended to handle non-linearity through techniques such as polynomial features, interaction terms, and kernel methods. Alternatively, inherently non-linear models like decision trees, random forests, or neural networks can be used for better performance on non-linear data.

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

**Answer:** Adaboost (Adaptive Boosting) and Gradient Boosting are two popular boosting algorithms used in machine learning to improve the performance of weak learners (usually decision trees) by combining them into a strong learner. Although both techniques are based on the principle of boosting, they differ in their approach to combining the weak learners and updating their weights. Here are the key differences between Adaboost and Gradient Boosting:

## Adaboost (Adaptive Boosting)

1. **Weight Update Mechanism:**
  - Adaboost assigns weights to each training sample and adjusts these weights after each iteration. Samples that are misclassified by the current weak learner have their weights increased, so that subsequent learners focus more on these hard-to-classify samples.
  - The weight update formula involves an exponential function of the prediction error.
2. **Error Minimization:**
  - Adaboost aims to minimize the exponential loss function. It adjusts the weights of the misclassified samples in such a way that subsequent weak learners are forced to correct the mistakes of the previous ones.
3. **Sequential Learning:**
  - Adaboost builds the model sequentially, where each new model is trained to correct the errors of the previous models by focusing on the hard-to-classify samples.
  - The final model is a weighted sum of the weak learners, where each learner's contribution is weighted according to its accuracy.
4. **Algorithm Simplicity:**
  - Adaboost is relatively simpler to implement compared to Gradient Boosting. It primarily involves updating the weights of the samples and combining the weak learners.
5. **Usage:**
  - Adaboost is commonly used with decision stumps (one-level decision trees) but can be used with any weak learner.
  - It is more sensitive to noisy data and outliers because it aggressively focuses on hard-to-classify samples.

## Gradient Boosting

1. **Weight Update Mechanism:**
  - Gradient Boosting does not directly adjust sample weights. Instead, it fits the new learner to the negative gradient of the loss function with respect to the current model predictions. This negative gradient is essentially the residual errors of the current model.
2. **Error Minimization:**
  - Gradient Boosting aims to minimize a specified loss function (e.g., mean squared error for regression, log-loss for classification) by iteratively adding weak learners that correct the residual errors of the combined ensemble.
  - It is more flexible in terms of the loss functions it can minimize.
3. **Sequential Learning:**
  - Like Adaboost, Gradient Boosting builds the model sequentially. However, each new model is trained to predict the residual errors (gradients) of the previous models.
  - The final model is a sum of the weak learners, each trained to correct the mistakes of the previous learners.
4. **Algorithm Complexity:**
  - Gradient Boosting is more complex and flexible than Adaboost. It involves choosing an appropriate loss function and optimizing it using gradient descent.
5. **Usage:**
  - Gradient Boosting can use any differentiable loss function, making it highly versatile. Common implementations include Gradient Boosting Machines (GBM), XGBoost, LightGBM, and CatBoost.

- It is more robust to overfitting and can handle noisy data better than Adaboost by incorporating regularization techniques.

**Comparison Table**

Feature	Adaboost	Gradient Boosting
<b>Weight Update</b>	Adjusts weights of misclassified samples	Fits learners to negative gradient of loss function
<b>Error Minimization</b>	Minimizes exponential loss function	Minimizes specified loss function
<b>Learning Process</b>	Sequential, focuses on misclassified samples	Sequential, fits learners to residual errors
<b>Model Combination</b>	Weighted sum of weak learners	Sum of weak learners
<b>Algorithm Simplicity</b>	Simpler	More complex and flexible
<b>Flexibility</b>	Less flexible, typically uses decision stumps	Highly flexible, supports various loss functions
<b>Sensitivity to Noise</b>	More sensitive	More robust with regularization
<b>Common Implementations</b>	Original Adaboost algorithm	GBM, XGBoost, LightGBM, CatBoost

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

**Answer:**

The bias-variance trade-off is a fundamental concept in machine learning that describes the trade-off between two sources of error that affect the performance of predictive models: bias and variance. Understanding this trade-off is crucial for building models that generalize well to new, unseen data.

### **Bias**

**Definition:** Bias refers to the error introduced by approximating a real-world problem, which may be complex, by a much simpler model. It is the difference between the average prediction of our model and the correct value we are trying to predict.

### **Characteristics:**

- **High Bias:** A model with high bias makes strong assumptions about the data, leading to systematic errors. High bias can cause underfitting, where the model is too simple to capture the underlying patterns in the data.
- **Low Bias:** A model with low bias makes fewer assumptions about the data and can capture more complex relationships.

### **Variance**

## Trade-Off

**Explanation:** The bias-variance trade-off describes the balance between bias and variance.

Increasing model complexity tends to decrease bias but increase variance, while decreasing model complexity tends to increase bias but decrease variance. The goal is to find the right balance where both bias and variance are minimized to achieve low overall error.

- **Underfitting:** High bias and low variance. The model is too simple to capture the underlying patterns (e.g., linear regression on a non-linear dataset).
- **Overfitting:** Low bias and high variance. The model captures noise in the training data, leading to poor generalization (e.g., a very deep decision tree).

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

### Answer:

In Support Vector Machines (SVMs), kernels are functions used to transform the input data into a higher-dimensional space where it is easier to perform the classification or regression tasks. Different types of kernels can be used depending on the nature of the data. Here are short descriptions of the Linear, Radial Basis Function (RBF), and Polynomial kernels commonly used in SVM:

#### Linear Kernel

**Description:** The linear kernel is the simplest kernel function. It is used when the data is linearly separable, meaning that the classes can be separated with a straight line (or hyperplane in higher dimensions).

#### Usage:

- Suitable for linearly separable data.
- Computationally efficient and simple to implement.
- Often used as a baseline for comparing with more complex kernels.

**Example:** Classifying linearly separable data like email spam detection or sentiment analysis where the data can be linearly separated in the feature space.

#### Radial Basis Function (RBF) Kernel

**Description:** The RBF kernel, also known as the Gaussian kernel, is a popular choice for non-linear data. It maps the data into an infinite-dimensional space, allowing the SVM to create a complex decision boundary.

#### Usage:

- Suitable for non-linear data where the relationship between the features and the target variable is complex.
- Flexible and capable of handling various types of data distributions.

**Example:** Image classification, handwriting recognition, and other tasks with non-linear relationships between features.

### **Polynomial Kernel**

**Description:** The polynomial kernel represents the similarity of vectors in a feature space over polynomials of the original variables, allowing the SVM to fit non-linear decision boundaries.

#### **Usage:**

- Suitable for data where the relationship between features and the target variable can be represented by polynomial functions.
- Can handle interactions between features.
- Requires tuning of the degree  $d$  and other parameters.

**Example:** Gene classification, text classification, and other tasks where polynomial relationships are expected.