

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 -In regression analysis, both R-squared and the Residual Sum of Squares (RSS) are goodness of fit of a model.

R-squared: R-squared is the proportion of the variance in the dependent variable that is predictable from the independent variables. As R^2 is equal to 1 (regression model perfectly fits the data) or 0 (indicates that the model does not explain any of the variability of the response data when its mean). Easy to interpret, goodness of fit with the same dependent variable and misleading when non-linear models or models with many predictors.

$$R^2 = 1 - (RSS/TSS)$$

- provides a normalized measure

Residual Sum of Squares (RSS): RSS is the sum of the squared differences between the observed values and the values predicted by the model. It is a measure of the total deviation of the response values from the fit to the response values.

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

- Lower RSS indicates a better fit of the model to the data.
- It provides an absolute measure of fit, reflecting the magnitude of prediction errors.
- It is not standardized, making it difficult to compare across different datasets or models with different scales.

R-squared is typically more convenient for comparing models and understanding the proportion of explained variance.

2. What are TSS (Total Sum of Squares), ESS (Explained Sum of Squares) and RSS (Residual Sum)?

Ans- **Total Sum of Squares (TSS):** TSS measures the total variation in the dependent variable. It is the sum of the squares of the differences between each observed value and the mean of the observed values.

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

where y is the observed value

\bar{y} is the mean of the dependent values.

Explained Sum of Squares (ESS): Measures the variation in the dependent variable that is explained by the independent variables in the model. It is the sum of the squares of the differences between the predicted values and the mean of the observed values.

$ESS = \sum (y_i - \bar{y})^2$ where \hat{y}_i is the i -th predicted value from the regression model.

Residual Sum of Squares (RSS): RSS is the sum of the squared differences between the observed values and the values predicted by the model. It is a measure of the total deviation of the response values from the fit to the response values.

$RSS = \sum (y_i - \hat{y}_i)^2$

- Lower RSS indicates a better fit of the model to the data.
- It provides an absolute measure of fit, reflecting the magnitude of prediction errors.
- It is not standardized, making it difficult to compare across different datasets or models with different scales.

The following equation shows the relationship:

$TSS = ESS + RSS$

- The Total Sum of Squares (TSS) represents the total variability in the dependent variable.
- The Explained Sum of Squares (ESS) represents the portion of the total variability that is explained by the regression model.
- The Residual Sum of Squares (RSS) represents the portion of the total variability that is not explained by the model.

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

Regularization is essential in machine learning to create models that generalize well to unseen data. By penalizing excessive complexity and preventing overfitting, regularization techniques contribute to more reliable and interpretable models, especially in the presence of high-dimensional data or multicollinearity.

1. Preventing Overfitting -occurs when underlying patterns in the training data but also the noise and random fluctuations. This leads to excellent performance on training data but poor generalization to new, unseen data.

2. Improving Generalization- model which is complex (degree is high) may fit the training data very well but fail to perform well on test data.

3. Handling Multicollinearity- Multicollinearity occurs when two or more predictor variables in a regression model are highly correlated. This can cause instability in the estimates of regression coefficients, making the model unreliable.

4. Feature Selection-In datasets with a large number of features, not all features are necessarily relevant. Including irrelevant or redundant features can harm the model's performance.

Types of Regularization

- **L1 Regularization (Lasso)** =Original Loss+ $\lambda \sum_j |\beta_j|$ \text{Loss}
- **L2 Regularization (Ridge)** =Original Loss+ $\lambda \sum_j \beta_j^2$ \text{Loss}
- **Elastic Net**: combines both L1 and L2 regularization penalties.

4.What is Gini–impurity index?

Ans-The Gini impurity index is a metric used to evaluate the purity of a node in decision tree algorithms, particularly in classification tasks. It measures how often a randomly chosen element from the set would be incorrectly labeled if it were randomly labeled according to the distribution of labels in the node. The Gini impurity index is used to decide the best split at each step of building a decision tree.

For a node t in a decision tree, let $p(i|t)$ be the proportion of instances belonging to class i in node t . The Gini impurity $G(t)$ is defined as:

$$G(t) = 1 - \sum_{i=1}^C p(i|t)^2$$

where C is the number of classes.

Interpretation

- **Range:** The Gini impurity ranges from 0 to $(1 - 1/C)$, where C is the number of classes. For a binary classification problem, this means the Gini impurity ranges from 0 to 0.5.
- **Purity:**
 1. $G(t) = 0$: The node is pure, meaning all instances belong to a single class.
 2. $G(t)$ close to 1: The node is impure, meaning the instances are evenly distributed among the different classes.

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

Ans- Yes, unregularized decision trees are prone to overfitting because they can create overly complex models that capture noise and specific patterns in the training data. Regularization techniques are essential to control the complexity of decision trees, thereby improving their ability to generalize to new, unseen data and enhancing overall model performance.

- **High Flexibility and Complexity**
- **2. Lack of Constraints**

- 3. Absence of Pruning
- Evidence of Overfitting

To reduce the risk of overfitting, various regularization techniques can be applied:

- Limit Tree Depth:
- Minimum Samples for Splits and Leaves
- Maximum Number of Leaf Nodes
- Pruning
- Using Ensemble Methods.

6. What is an ensemble technique in machine learning?

Ensemble techniques in machine learning involve combining the predictions of multiple models to produce a more accurate and robust prediction than any single model could achieve on its own. The primary goal of ensemble methods is to improve the generalization ability of models by reducing variance, bias, or improving predictions. Here are some common ensemble techniques:

- Bagging (Bootstrap Aggregating)
- Boosting
- Stacking (Stacked Generalization)
- Voting

Advantages of Ensemble Techniques

- Improved Accuracy
- Reduced Overfitting
- Robustness

Disadvantages of Ensemble Techniques

- Increased Complexity
- Interpretability

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

Bagging and Boosting are powerful ensemble methods with distinct strategies for improving model performance. Bagging aims to reduce variance by averaging over multiple models trained in parallel on different subsets of the data, while Boosting aims to reduce both bias and variance by training models sequentially, with each model focusing on the errors of its predecessors. The choice between bagging and boosting depends on the specific problem and the desired balance between bias and variance reduction.

Aspect	Bagging	Boosting
Training Approach	Parallel	Sequential
Data Handling	Bootstrap sampling (with replacement)	Adjusting weights of samples
Focus	Reducing variance	Reducing bias and variance
Model Independence	Models trained independently	Models trained dependently
Model Complexity	Often uses high-variance models	Often uses weak learners
Combination Method	Majority vote (classification) or average (regression)	Weighted sum of models
Example	Random Forest	AdaBoost, Gradient Boosting

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

Out-of-bag (OOB) error is an important concept in the context of Random Forests, an ensemble learning method. It provides a way to estimate the prediction error of the model without the need for a separate validation or test dataset. Here's a detailed explanation:

Out-of-Bag Error in Random Forests

Bootstrap Sampling: Random Forests, each decision tree is trained on a different bootstrap sample of the training data. A bootstrap sample is created by randomly selecting samples from the original dataset with replacement, meaning some samples may appear multiple times, while others may not appear at all.

1. **Out-of-Bag (OOB) Samples:** For each bootstrap sample used to train a tree, there will be some data points from the original dataset that are not included in that sample. These data points are known as out-of-bag samples.
2. **Estimating OOB Error:** After a tree is trained, it can be tested on its corresponding OOB samples. The prediction error for these OOB samples is calculated and recorded.
3. This process is repeated for each tree in the Random Forest, resulting in multiple OOB error estimates.
4. **Aggregating OOB Error:** The overall OOB error for the Random Forest is calculated by aggregating the OOB error estimates from all trees. Typically, this is done by averaging the prediction errors of the OOB samples across all trees.

Benefits of OOB Error

1. **Validation Without a Separate Dataset**
2. **Efficiency**
3. **Reliability**

The steps to calculate-

Train Each Tree -Train each decision tree on a bootstrap sample of the original data.

Predict OOB Samples: For each tree, predict the outcomes for the data points that were not included in its bootstrap sample (the OOB samples).

Compute OOB Error for Each Tree: Compute the error (e.g., misclassification error for classification, mean squared error for regression) for the predictions made on the OOB samples for each tree.

Aggregate OOB Errors:Average the errors across all trees to obtain the overall OOB error estimate for the Random Forest.

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

K-fold cross-validation is a robust and widely used technique for assessing the performance and generalization ability of a machine learning model. It involves partitioning the dataset into KKK equally sized (or nearly equally sized) subsets or folds and then performing the following steps:

Steps in K-Fold Cross-Validation

1. **Partition the Data:** Split the entire dataset into KKK folds. Each fold is a subset of the data, and the number KKK is a user-specified parameter (e.g., 5, 10).
2. **Training and Validation:**For each fold iii (where iii ranges from 1 to KKK):
 - **Training Set:** Use $K-1$ folds combined as the training set.
 - **Validation Set:** Use the remaining fold as the validation set.
 - Train the model on the training set and evaluate it on the validation set.
3. **Repeat:**Repeat the process KKK times, each time with a different fold serving as the validation set and the remaining $K-1$ folds as the training set.
4. **Average Performance:**Calculate the performance metric (e.g., accuracy, mean squared error) for each fold.
5. Average the performance metrics across all KKK folds to obtain the overall performance estimate.

Advantages of K-Fold Cross-Validation

1. Better Utilization of Data
2. Reduced Variance
3. Model Selection

Disadvantages-

Computationally Intensive and I imbalanced Data

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

Hyperparameter tuning in machine learning refers to the process of finding the best set of hyperparameters for a model. Hyperparameters are parameters that are set before the learning process begins and control the learning process itself. Unlike model parameters, which are learned during training (e.g., weights in neural networks), hyperparameters are typically specified by the practitioner based on domain knowledge, trial and error, or best practices.

It involves experimenting with different hyperparameter configurations and selecting the best combination that maximizes the model's effectiveness and efficiency.

1. Optimizing Model Performance
2. Improving Generalization
3. Tailoring to the Problem
4. Enhancing Efficiency
5. Model Interpretability and Stability
6. Hyperparameter tuning aligns

Techniques for Hyperparameter Tuning

1. Grid Search
2. Random Search
3. Bayesian Optimization
4. Gradient-based Optimization
5. Automated Hyperparameter Tuning Tools

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

Large learning rate in Gradient Descent can lead to several issues, primarily affecting the convergence and stability of the optimization process. Main issues are-

1. Overshooting the Minimum
2. Divergence

- 3. Slow Convergence
- 4. Poor Generalization
- 5. Instability in Training

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

Logistic Regression is a linear model designed primarily for binary classification tasks, where it models the probability of a binary outcome based on linear predictors. The decision boundary in logistic regression is a linear function of the input features, separating the classes with a straight line or hyperplane in higher dimensions.

Linear Decision Boundary: Non-linearly separable data, where classes are not separable by a straight line or hyperplane in the input space, logistic regression will not perform well.

Underfitting: When applied to non-linear data, logistic regression may underfit because it cannot capture the complex relationships and interactions between features that non-linear models can. This limitation can lead to poor performance and low accuracy on datasets with non-linear class boundaries.

Feature Transformation- While logistic regression itself cannot model non-linear relationships. However, this approach has its limitations and may not adequately capture highly complex non-linear relationships.

To handle non-linear data effectively - **Decision Trees and Ensemble Methods**, SVM's, Neural Networks and kernel logistic Regression can be used.

13. Differentiate between Adaboost and Gradient Boosting.

Adaboost (Adaptive Boosting) and Gradient Boosting are both ensemble learning techniques that sequentially combine weak learners (typically decision trees) to create a strong predictive model. Despite some similarities, they differ significantly in their approach to boosting and how they update the model in each iteration.

Adaboost (Adaptive Boosting)

1. **Approach:** Adaboost focuses on adjusting the weights of training instances based on their classification error. It assigns higher weights to incorrectly classified instances so that subsequent weak learners focus more on these difficult cases.

2. **Weighted Training:** Each base learner (weak learner), often a decision tree with limited depth (stump), is trained sequentially. After each iteration, the weights of misclassified data points are increased.
3. **Combining Predictions:** Predictions from each weak learner are combined using a weighted majority vote (for classification) or a weighted sum (for regression).
4. **Advantages:** Adaboost is straightforward to implement and can achieve good results with relatively few hyperparameters to tune.
5. It can handle complex interactions in the data well and tends to be less prone to overfitting.
6. **Disadvantages:** Adaboost can be sensitive to noisy data and outliers because it tends to focus more on difficult instances.
7. It can suffer from overfitting if the base learner is too complex or if the number of iterations (weak learners) is too high.

Gradient Boosting

1. **Approach:** Gradient Boosting builds an ensemble of trees in a sequential manner, where each tree is trained to correct the errors (residuals) of the previous tree.
2. **Residual Fitting:** Unlike Adaboost, which adjusts weights of instances, Gradient Boosting fits each new tree to the residual errors made by the previous trees.
3. **Gradient Descent Optimization:** The learning process in Gradient Boosting involves minimizing a loss function (e.g., mean squared error for regression, cross-entropy loss for classification) by gradient descent. Each new tree is trained to reduce the residual error of the ensemble.
4. **Combining Predictions:** Predictions from all trees (weak learners) are summed to make the final prediction. Gradient Boosting usually uses a shrinkage parameter (learning rate) to control the contribution of each tree.
5. **Advantages:** Gradient Boosting tends to yield better performance than Adaboost on most datasets, as it can capture complex relationships in the data through iterative refinement. It can handle both regression and classification tasks effectively.
6. **Disadvantages:** Gradient Boosting can be more complex to implement and tune compared to Adaboost, requiring careful tuning of hyperparameters such as learning rate, tree depth, and number of iterations. It may be more prone to overfitting if the model is too complex or if the learning rate is set too high.

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

The balance between the Bias error and the Variance error is the **Bias-Variance Tradeoff**. Helps in understanding the overall performance and generalization ability of a model.

Bias refers to the error introduced by approximating a real-world problem with a simplified model. Bias

refers to the error introduced by approximating a real-world problem with a simplified model

Variance refers to the amount by which the model's predictions vary for different training datasets. It measures the sensitivity of the model to the randomness in the training data. High variance can lead to overfitting, where the model performs very well on the training data but poorly on new, unseen data. The model memorizes the noise and specific details of the training data rather than generalizing from the underlying patterns.

Trade-off

- **Goal:** The goal in machine learning is to find a balance between bias and variance to achieve a model that generalizes well to new data.
- **Challenges:** Decreasing bias typically increases variance, and vice versa. Therefore, reducing one type of error often increases the other, leading to a trade-off.
- **Optimal Point:** The optimal model complexity is typically found at the point where the total error (sum of bias and variance) is minimized.
- **Regularization:** Techniques like regularization can help control model complexity and reduce variance, thereby improving generalization.
- **Cross-validation:** Cross-validation techniques like K-fold cross-validation help in assessing the bias-variance trade-off by providing estimates of both bias and variance.

Practical Considerations

- **Model Selection:** Choosing the right model architecture and complexity is crucial.
- **Data Size:** Larger datasets can help reduce variance by providing more examples for the model to learn from, potentially reducing overfitting.
- **Hyperparameter Tuning:** Adjusting hyperparameters (e.g., learning rate, regularization strength) can influence the bias-variance trade-off. For instance, increasing regularization strength tends to reduce variance but may increase bias.

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

1. Linear Kernel

The Linear kernel is the simplest and efficient of kernel. Data is linearly separable as data is classified into 2 classes using a straight line.(2D)

$$K(x, x') = x^T x'$$

2. RBF (Radial Basis Function) Kernel—

Versatile and effective for capturing complex non-linear relationships, controlled by γ . It can handle non-linear decision boundaries and is effective for datasets with complex relationships. It

measures the similarity between two data points based on the Euclidean distance in the transformed feature space.

$$K(x, x') = \exp(-\gamma \|x - x'\|^2)$$

3. Polynomial Kernel

The Polynomial kernel computes the similarity between two vectors. It is used when data points may be separable by a polynomial decision boundary in space (polynomials), d is the degree of the polynomial, and c is an optional coefficient that can be adjusted.

$$K(x, x') = (x^T x' + c)^d$$