# MACHINE LEARNING

# ASSIGNMENT - 5

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

1. R-squared or Residual Sum of Squares (RSS) which one of these two is a better measure of

goodness of fit model in regression and why?

Answer: R-squared (coefficient of determination) is generally considered a better measure of the goodness of fit in regression models compared to the Residual Sum of Squares (RSS). Here's why:

Interpretability:

R-squared provides a more interpretable measure of goodness of fit. It represents the proportion of the variance in the dependent variable that is explained by the independent variables in the model. Higher R-squared values indicate a better fit.

RSS, on the other hand, is an absolute measure of the total sum of squared residuals, but it doesn't offer a clear interpretation of the model's explanatory power.

Normalization:

R-squared is normalized, meaning it takes values between 0 and 1. A value of 1 indicates a perfect fit, where the model explains all the variability in the dependent variable. A value of 0 indicates that the model does not explain any variability.

RSS is not normalized, and its interpretation depends on the scale of the dependent variable. It doesn't provide a clear indication of the proportion of variability explained.

Comparison Across Models:

R-squared allows for easy comparison between different models. When comparing models, a higher R-squared suggests a better fit.

RSS alone doesn't provide a basis for comparing models unless you consider other factors like the number of variables. It doesn't account for the overall variability in the dependent variable.

Adjustment for Degrees of Freedom:

Adjusted R-squared further adjusts for the number of independent variables in the model, penalizing the inclusion of unnecessary variables. This is important in avoiding overfitting.

RSS doesn't have a built-in adjustment for the number of variables, which means it might not penalize overfitting.

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.

1.  Answer: **Total Sum of Squares (TSS):**
    *   TSS represents the total variability in the dependent variable (y).

- It measures the total deviation of each data point from the mean of the dependent variable.
- Equation: $TSS = \sum(y_i - \bar{y})^2$, where $\bar{y}$ is the mean of the dependent variable.

2. **Explained Sum of Squares (ESS):**
   - ESS measures the variability in the dependent variable that is explained by the regression model.
   - It is the sum of squared differences between the predicted values (ŷ) and the mean of the dependent variable.
   - Equation: $ESS = \sum(\hat{y}_i - \bar{y})^2$, where $\hat{y}_i$ is the predicted value for each observation.

3. **Residual Sum of Squares (RSS):**
   - RSS measures the unexplained variability in the dependent variable, which is attributed to the residuals (the differences between the observed and predicted values).
   - It is the sum of squared residuals.
   - Equation: $RSS = \sum(y_i - \hat{y}_i)^2$, where $y_i$ is the observed value for each observation.

4. **Equation Relating TSS, ESS, and RSS:**
   - The relationship between TSS, ESS, and RSS is given by the fundamental identity: $TSS = ESS + RSS$
   - This identity states that the total variability in the dependent variable (TSS) can be decomposed into the variability explained by the model (ESS) and the unexplained variability (RSS).


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

Regularization is a technique used in machine learning to prevent overfitting and improve the generalization performance of a model. The primary needs and motivations for regularization in machine learning are:

**Preventing Overfitting:**

Overfitting occurs when a model learns the training data too well, capturing noise and irrelevant details. As a result, the model performs poorly on new, unseen data.

Regularization helps prevent overfitting by adding a penalty term to the model's objective function. This penalty discourages the model from assigning excessively large weights to features, making the model more robust to noise in the training data.

**Handling Multicollinearity:**

Multicollinearity arises when two or more features in a dataset are highly correlated. This can lead to instability in the model, making it sensitive to small changes in the input data.

Regularization techniques, such as Ridge Regression, can mitigate the impact of multicollinearity by imposing constraints on the magnitudes of the coefficients.

**Feature Selection:**

Regularization methods can encourage sparsity in the model by driving some feature weights to exactly zero. This leads to automatic feature selection, where only the most informative features are retained, and irrelevant features are effectively ignored.

Lasso Regression is a regularization technique known for inducing sparsity in the model.

**Improving Model Interpretability:**

Regularization can simplify the model and make it more interpretable by reducing the number of parameters. A simpler model is often easier to understand and explain.

In scenarios where model interpretability is important, regularization provides a balance between complexity and performance.

**Handling Noisy Data:**

Noisy data points can have a disproportionate impact on the model's performance. Regularization helps the model generalize better by downweighting the influence of noisy data points during training.

It adds a penalty for complex models, discouraging the model from fitting the noise in the training data.

**Balancing Bias and Variance:**

Regularization helps strike a balance between bias and variance in the model. A model with too much complexity may have low bias but high variance, leading to overfitting. Regularization helps control the trade-off between bias and variance.

4.What is Gini–impurity index?

The Gini impurity index, often referred to as Gini index or Gini impurity, is a measure of the impurity or disorder in a set of categorical data. It is commonly used in decision tree algorithms for binary classification problems. The Gini index quantifies the likelihood of an incorrect classification of a randomly chosen element in the dataset.

For a binary classification problem with two classes (usually denoted as 0 and 1), the Gini impurity index for a node in a decision tree is calculated using the following formula:

$Gini(D) = 1 - \sum_{i=1}^{c} p_i^2$

where:

$D$ is the dataset associated with a particular node.

$c$ is the number of classes.

$p_i$ is the proportion of instances of class $i$ in the node.

The Gini index ranges from 0 to 1, where:

$Gini(D) = 0$ indicates perfect purity (all instances belong to the same class).

$Gini(D) = 1$ indicates maximum impurity (an equal distribution of instances across all classes).

In the context of a decision tree, the Gini index is used to evaluate the quality of a split. When deciding which feature to split on and what threshold to use, the algorithm chooses the split that minimizes the Gini index in the resulting child nodes. The split that reduces impurity the most is considered the best.

The Gini impurity index is an alternative to the entropy-based information gain criterion and is widely used in decision tree algorithms, including CART (Classification and Regression Trees). Both

Gini impurity and information gain serve as criteria for creating decision tree nodes that effectively partition the data based on feature values.

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

 **Answer**: Yes, unregularized decision trees are prone to overfitting, and the primary reason for this vulnerability lies in their inherent nature of being able to create highly complex and detailed structures to fit the training data. Overfitting occurs when a model captures noise and specific details in the training data to an extent that it negatively impacts its ability to generalize to new, unseen data. Here's why unregularized decision trees are susceptible to overfitting:

1. **Flexibility to Capture Complex Patterns**
2. **Memorization of Training Data**
3. **Sensitive to small variations**
4. **Large and deep trees**
5. **High variance**

6.  What is an ensemble technique in machine learning?

Answer: Ensemble learning is a generic approach to machine learning that aims to improve predictive performance by mixing predictions from several models.
Despite the fact that you can create an almost infinite number of ensembles for any predictive modeling problem, there are three strategies that dominate the field of ensemble machine learning. So much so that instead of algorithms as such, each is a subject of study that has given rise to other more specialized approaches. Using ensemble approaches, a single best-fit prediction model can be generated by combining multiple basic models.
Bagging, stacking, and boosting are the three basic types of ensemble learning techniques, and understanding each approach in depth and incorporating it into your predictive modeling project is critical

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

| S.NO | Bagging | Boosting |
|------|---------|----------|
| 1. | The simplest way of combining predictions that belong to the same type. | A way of combining predictions that belong to the different types. |
| 2. | Aim to decrease variance, not bias. | Aim to decrease bias, not variance. |
| 3. | Each model receives equal weight. | Models are weighted according to their performance. |
| 4. | Each model is built independently. | New models are influenced by the performance of previously built models. |

| S.NO | Bagging | Boosting |
|---|---|---|
| 5. | Different training data subsets are selected using row sampling with replacement and random sampling methods from the entire training dataset. | Every new subset contains the elements that were misclassified by previous models. |
| 6. | Bagging tries to solve the over-fitting problem. | Boosting tries to reduce bias. |
| 7. | If the classifier is unstable (high variance), then apply bagging. | If the classifier is stable and simple (high bias) the apply boosting. |
| 8. | In this base classifiers are trained parallelly. | In this base classifiers are trained sequentially. |
| 9 | Example: The Random forest model uses Bagging. | Example: The AdaBoost uses Boosting techniques |

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

 **Answer**: A random forest is an ensemble machine-learning model that is composed of multiple decision trees. A decision tree is a model that makes predictions by learning a series of simple decision rules based on the features of the data. A random forest combines the predictions of multiple decision trees to make more accurate and robust predictions.

Random Forests are often used for classification and regression tasks. In classification, the goal is to predict the class label (e.g., "cat" or "dog") of each sample in the dataset. In regression, the goal is to predict a continuous target variable (e.g., the price of a house) based on the features of the data.
Random forests are popular because they are easy to train, can handle high-dimensional data, and are highly accurate. They also have the ability to handle missing values and can handle imbalanced datasets, where some classes are more prevalent than others.

To train a random forest, you need to specify the number of decision trees to use (the n_estimators parameter) and the maximum depth of each tree (the max_depth parameter). Other hyperparameters, such as the minimum number of samples required to split a node and the minimum number of samples required at a leaf node, can also be specified.
Once the random forest is trained, you can use it to make predictions on new data. To make a prediction, the random forest uses the predictions of the individual decision trees and combines them using a majority vote or an averaging technique.

9 . What is K-fold cross-validation?

**Answer:** K-fold cross-validation is a resampling technique used in machine learning for model evaluation and hyperparameter tuning. It involves partitioning the original training data set into K subsets (folds) of approximately equal size. The learning algorithm is then trained K times, each time using K-1 folds for training and the remaining fold for validation or testing. This process is repeated K times, with each of the K folds used exactly once as a validation data set. The main steps of K-fold cross-validation are as follows:

Partitioning:

The original dataset is randomly divided into K equally sized folds.

Each fold is treated as a testing set, and the remaining K-1 folds are used for training.

Training and Evaluation:

The learning algorithm is trained K times, each time on a different combination of K-1 folds.

After training, the model is evaluated on the fold that was not used for training.

Performance Metrics:

Performance metrics, such as accuracy, precision, recall, or mean squared error, are computed for each iteration.

Average Performance:

The performance metrics obtained from each iteration are averaged to obtain a single performance estimate.

K-fold cross-validation provides several benefits:

Robustness: By repeating the training and testing process K times with different subsets, the evaluation process becomes more robust, and the performance estimate is less sensitive to the randomness of a single data split.

Better Utilization of Data: It allows the algorithm to be trained and tested on different portions of the dataset, ensuring that the model sees and learns from the entire dataset.

Model Variance Evaluation: K-fold cross-validation helps in assessing the variability of the model's performance by providing K different performance estimates.

Common choices for the value of K are 5 and 10, but the optimal choice depends on factors such as the size of the dataset and the computational resources available. In cases where the dataset is small, leave-one-out cross-validation (LOOCV), a special case of K-fold cross-validation with K equal to the number of instances, may be used.

K-fold cross-validation is a valuable tool for model assessment and hyperparameter tuning, providing a more reliable estimate of a model's performance than a single train/test split.


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

Answer: Hyperparameter tuning is the process of selecting the optimal values for a machine learning model's hyperparameters. Hyperparameters are settings that control the learning process of the model, such as the learning rate, the number of neurons in a neural network, or the kernel size in a support vector machine. The goal of hyperparameter tuning is to find the values that lead to the best performance on a given task.

n the context of machine learning, hyperparameters are configuration variables that are set before the training process of a model begins. They control the learning process itself, rather than being learned from the data. Hyperparameters are often used to tune the performance of a model, and they can have a significant impact on the model's accuracy, generalization, and other metrics.

The purpose of hyperparameter tuning is to find the best set of hyperparameters for a given machine learning model. This can improve the model's performance on unseen data, prevent overfitting, and reduce training time.

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

Answer: Having a large learning rate in gradient descent can lead to several issues, and it can negatively impact the training process of a machine learning model. Here are some common issues associated with a large learning rate:

Divergence:

With a large learning rate, the updates to the model parameters may become too large, causing the optimization process to diverge rather than converge.

Divergence means that the parameter values oscillate or move away from the optimal values, preventing the model from reaching a minimum.

Overshooting the Minimum:

Large learning rates can cause the optimization algorithm to overshoot the minimum of the cost function.

Instead of converging to the minimum, the algorithm may oscillate around it or even move away, preventing the model from reaching optimal parameter values.

Unstable Training:

The training process becomes unstable, and the model's performance may fluctuate significantly between iterations.

The model may fail to converge to a solution, making it challenging to determine the best set of parameters.

Inability to Converge:

The optimization process may fail to converge to a minimum, and the algorithm may not reach a stable solution.

This can result in the model not learning from the data effectively and producing poor generalization performance.

NaN or Overflow Issues:

Large learning rates can lead to numerical instability, causing intermediate values in the optimization process to become extremely large (overflow) or result in NaN (Not a Number) values.

These numerical issues can halt the training process.

To address these issues, it's essential to choose an appropriate learning rate for the specific optimization problem. Techniques such as learning rate scheduling, adaptive learning rates (e.g., AdaGrad, RMSprop, Adam), and careful tuning can help prevent the problems associated with a large learning rate. Cross-validation and monitoring the training/validation performance during the training process can also guide the selection of an appropriate learning rate.

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

Answer: Logistic Regression is a linear model designed for binary classification tasks, and it makes predictions using a linear combination of input features. Despite its name, logistic regression is a classification algorithm, not a regression algorithm. It models the probability that an instance belongs to a particular class and uses the logistic (sigmoid) function to squash the linear prediction into the range [0, 1].

While Logistic Regression is inherently a linear model, it can be used for some non-linear decision boundaries, particularly when employing techniques like feature engineering or polynomial features. However, its ability to handle highly non-linear data is limited compared to more complex non-linear models, such as decision trees, support vector machines with non-linear kernels, or neural networks.

Here are some points to consider:

Linear Decision Boundaries:

Logistic Regression assumes a linear decision boundary in the feature space. If the relationship between features and the target variable is genuinely non-linear, Logistic Regression may struggle to capture complex patterns.

Feature Engineering:

Feature engineering can enhance Logistic Regression's ability to model non-linear relationships. Techniques like adding interaction terms, polynomial features, or transforming existing features can help the model capture non-linearities.

Limited Capacity:

Logistic Regression has limited capacity to learn complex patterns compared to more flexible models. For highly non-linear datasets, models that can capture intricate structures in the data may provide better performance.

Other Non-Linear Models:

For inherently non-linear problems, it is often more effective to use models explicitly designed to handle non-linearities, such as decision trees, kernelized support vector machines, or deep neural networks.

In summary, while Logistic Regression can be used for classification tasks with some non-linear aspects, it may not be the best choice for highly non-linear data. For such cases, more complex models capable of capturing intricate relationships in the data might offer better performance. It's essential to evaluate different models and choose the one that best fits the nature of the data and the complexity of the underlying patterns.

13. Differentiate between Adaboost and Gradient Boosting.

Answer:

| S.No | Adaboost | Gradient Boost |
|---|---|---|
| 1 | An additive model where shortcomings of previous models are identified by high-weight data points. | An additive model where shortcomings of previous models are identified by the gradient. |

| 2 | The trees are usually grown as decision stumps. | The trees are grown to a greater depth usually ranging from 8 to 32 terminal nodes. |
|---|---|---|
| 3 | Each classifier has different weights assigned to the final prediction based on its performance. | All classifiers are weighed equally and their predictive capacity is restricted with learning rate to increase accuracy. |
| 4 | It gives weights to both classifiers and observations thus capturing maximum variance within data. | It builds trees on previous classifier's residuals thus capturing variance in data. |

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

Answer: The bias-variance tradeoff is a fundamental concept in machine learning and statistics. It refers to the delicate balance between two sources of error in a predictive model: bias and variance. Bias represents the error due to overly simplistic assumptions in the learning algorithm. High bias can cause the model to underfit the data, leading to poor performance on both training and unseen data. Variance, on the other hand, reflects the model's sensitivity to small fluctuations in the training data. High variance can lead to overfitting, where the model captures noise in the training data and performs poorly on new, unseen data. The goal is to find the right level of complexity in a model to minimize both bias and variance, achieving good generalization to new data. Balancing these factors is essential for building models that perform well on a variety of datasets.

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

Answer: Support Vector Machines (SVMs) are a type of supervised learning model used for classification and regression tasks. SVMs use kernels to transform the input data into a higher-dimensional space, making it easier to find a hyperplane that separates different classes or captures complex relationships. Here's a short description of three commonly used kernels in SVM:

Linear Kernel:

Description: The linear kernel is the simplest and most straightforward SVM kernel. It represents the dot product of the input features in the original feature space.

Use Case: The linear kernel is suitable for linearly separable datasets or situations where a linear decision boundary is appropriate. It is computationally efficient and often used as a baseline.

Radial Basis Function (RBF) Kernel:

Description: The RBF kernel (also known as Gaussian kernel) measures the similarity between data points in a higher-dimensional space. It uses a radial basis function to transform the input features.

Use Case: The RBF kernel is versatile and effective for capturing complex, non-linear relationships in the data. It is widely used in SVMs and can handle cases where the decision boundary is not a straight line.

Polynomial Kernel:

Description: The polynomial kernel raises the dot product of the input features to a specified power, introducing non-linearity. It allows the SVM to capture polynomial relationships between features.

Use Case: The polynomial kernel is useful when the decision boundary is expected to be polynomial. It can model complex decision boundaries but may require tuning of the polynomial degree parameter.

Each kernel has its strengths and weaknesses, and the choice depends on the nature of the data and the underlying patterns. The RBF kernel is often a good default choice when the data is not inherently linear, but experimenting with different kernels and their parameters is recommended to find the best-performing model for a specific task.