# Machine Learning Assignment

1. R-squared is generally considered a better measure of goodness of fit for several reasons. However, RSS can be useful in specific contexts, such as when comparing models on the same dataset, particularly when the focus is on the actual errors rather than the proportion of explained variance.

2. TSS represents the total variability in the response variable.

   ESS represents the portion of the total variability that is explained by the model.

   RSS represents the portion of the total variability that is not explained by the model.

   Relationship Between TSS, ESS, and RSS

   These three metrics are related through the following equation:

   TSS=ESS+RSS

   This equation states that the total variance in the response variable (TSS) is equal to the sum of the variance explained by the model (ESS) and the variance that is not explained by the model (RSS).

3. Regularization is essential in machine learning for preventing overfitting, enhancing generalization, facilitating feature selection, and improving the stability of models. By adding a penalty for complexity, regularization techniques encourage the development of simpler, more robust models that perform better on unseen data.

4. The Gini impurity index is a measure used in decision trees and other machine learning algorithms to evaluate the purity of a dataset, particularly in classification tasks. It quantifies 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 subset.

5. Yes, unregularized decision trees are prone to overfitting. This is primarily because of their nature and how they function during the training process. Here are the main reasons why unregularized decision trees tend to overfit:

   1. Complexity and Depth
   2. High Variance
   3. Overfitting Indicators

6. An ensemble technique in machine learning refers to methods that combine multiple individual models to produce a single, more powerful predictive model. The key idea behind ensemble techniques is that combining multiple models can often lead to better performance than any single model, primarily by reducing errors due to bias, variance, or both.
   Techniques are : Bagging, boosting and stacking

7. Bagging and boosting are two popular ensemble techniques in machine learning that aim to improve the performance and robustness of models by combining multiple individual models. Despite their common goal, they achieve this in fundamentally different ways.

   Training Method:

   - Bagging: Models are trained independently and in parallel on different subsets of the data.
   - Boosting: Models are trained sequentially, with each model correcting the errors of the previous ones.

   Focus:

   - Bagging: Aims to reduce variance and improve stability.
   - Boosting: Aims to reduce bias and improve accuracy by focusing on hard-to-predict instances.

   Model Combination:

   - Bagging: Uses averaging (regression) or majority voting (classification) to combine models.
   - Boosting: Combines models through a weighted sum, giving more influence to better-performing models.

   Risk of Overfitting:

   - Bagging: Less prone to overfitting due to model independence.
   - Boosting: More prone to overfitting, especially if models are too complex or the boosting process is not properly regularized.

8. The out-of-bag (OOB) error is an important concept in the context of random forests, which is a bagging ensemble method primarily used for classification and regression tasks. OOB error provides an internal, unbiased estimate of the model's prediction error without the need for a separate validation set.

9. K-fold cross-validation is a widely used technique for assessing the performance and generalizability of a machine learning model. It involves partitioning the dataset into $KKK$ equal-sized subsets or folds and then performing the training and validation process $KKK$ times.

10. Hyperparameter tuning in machine learning is the process of optimizing the hyperparameters of a model to improve its performance. Hyperparameters are settings or configurations that are set before the learning process begins and are not learned from the data. They control the behavior of the training algorithm and the structure of the model.

11. Having a large learning rate in gradient descent can lead to several issues that can hinder the training process and the convergence of the model. Here are the main issues:

- Divergence:

  If the learning rate is too large, the updates to the model parameters (weights) can be too drastic. This can cause the parameters to oscillate around the minimum or even diverge away from it, preventing the model from converging to a good solution.

- Overshooting the Minimum:

  Large learning rates can cause the updates to overshoot the optimal values of the parameters. Instead of gradually descending towards the minimum of the loss function, the updates may move the parameters too far in each iteration, causing them to bounce back and forth across the minimum without settling down.

- Instability and Unpredictability:

  The training process becomes unstable and unpredictable with a large learning rate. Small changes in the training data or in the initialization of parameters can lead to significantly different outcomes, as the updates are overly sensitive to each data point.

- Failure to Converge:

  The gradient descent algorithm may fail to converge to the minimum of the loss function with a large learning rate. Instead of gradually reducing the loss, the updates might cause the loss to oscillate or even increase over time, preventing the algorithm from reaching a stable solution.

- Poor Generalization:

  Models trained with large learning rates are more likely to overfit the training data. This is because the rapid changes in parameter values can lead to a model that memorizes the training examples rather than learning general patterns that apply to unseen data.

- Computational Efficiency:

  While it may seem counterintuitive, very large learning rates can sometimes reduce the efficiency of gradient descent. This is because the algorithm may require more iterations to find a suitable solution due to the oscillations and instability caused by large updates.

12. Logistic Regression is a linear classification model, which means it assumes a linear relationship between the input variables (features) and the log-odds of the output (probability of a binary outcome). Here's why Logistic Regression is typically not suitable for handling non-linear data:

    Linear Decision Boundary
    - Linear Assumption:

      Logistic Regression assumes that the decision boundary separating the classes is linear. This means it can only learn to classify data points using a straight line (in 2D), a plane (in 3D), or a hyperplane (in higher dimensions).

    - Limitation on Complexity:

      If the relationship between the input features and the target variable is non-linear, Logistic Regression will struggle to capture and model that relationship accurately. It cannot learn more complex decision boundaries that are not linear.

13. Adaboost and Gradient Boosting are both popular ensemble learning techniques used to improve the performance of machine learning models. Despite sharing a goal of boosting weak learners (models that are slightly better than random guessing) into strong ones, they differ significantly in their approach and methodology

14. The bias-variance trade-off is a fundamental concept in supervised machine learning that helps us understand the sources of error in a model and how to balance them for optimal performance. It relates to the trade-off between a model's ability to capture underlying patterns (bias) and its sensitivity to noise or variability in the training data (variance).

15. Linear Kernel : The linear kernel is the simplest kernel function used in SVMs. It computes the dot product between the feature vectors in the original space.
    Radial Basis Function (RBF) Kernel : The RBF kernel, also known as the Gaussian kernel, maps the data into an infinite-dimensional space by computing the similarity between data points relative to a landmark or center.
    Polynomial Kernel : The polynomial kernel maps the data into a higher-dimensional space using polynomial functions of the original features

# Statistic Assignment Answer

1. D
2. C
3. C
4. B
5. C
6. B
7. A
8. A
9. B
10. A