# MACHINE LEARNING

R-squared (R²) is generally considered a better measure of goodness of fit in regression compared to Residual Sum of Squares (RSS).

1. **Interpretability**: R-squared provides a more intuitive interpretation. It represents the proportion of the variance in the dependent variable that is explained by the independent variables in the model. Higher values of R-squared indicate a better fit, with 1 indicating a perfect fit.

2. **Normalization**: R-squared is normalized, meaning its value is always between 0 and 1. This allows for easier comparison across different models and datasets. In contrast, RSS depends on the scale of the dependent variable, making it difficult to compare across models or datasets.

3. **Adjustment for Model Complexity**: R-squared adjusts for the number of predictors in the model through its adjusted version, Adjusted R-squared. This adjustment penalizes the addition of unnecessary variables that do not improve the model's fit. RSS, on the other hand, does not provide such adjustment and may favour more complex models even if they do not significantly improve the fit.

3.What is the need of regularization in machine learning?
1.Preventing Overfitting: Regularization helps prevent overfitting by imposing constraints on the model's parameters, discouraging them from taking on extreme or overly complex values that fit the training data too closely. By limiting the model's complexity, regularization encourages the learning of simpler and more generalizable patterns that are likely to perform well on unseen data.
2.Handling Multicollinearity: In linear regression and related models, multicollinearity occurs when predictor variables are highly correlated with each other. This can lead to unstable parameter estimates and inflated standard errors. Regularization techniques, such as Ridge regression, can mitigate multicollinearity by shrinking the coefficients of correlated predictors towards each other, improving the stability of the model.
3.Improving Model Interpretability: Regularization can lead to more interpretable models by discouraging the inclusion of irrelevant features or overly complex interactions. By penalizing unnecessary complexity, regularization encourages feature selection and model simplification, making it easier to understand the relationship between predictors and the target variable.

4. What is Gini–impurity index?
The Gini impurity index, often referred to as Gini impurity, is a measure of the impurity or uncertainty in a set of classification labels. It is commonly used as a criterion in decision tree algorithms for evaluating the quality of a split in a dataset.

In the context of binary classification problems, where there are two possible classes (e.g., positive and negative), the Gini impurity index for a set of labels is calculated using the following formula:
Gini impurity$=1-(p_1^2+p_2^2)$

5. <mark>Are unregularized decision-trees prone to overfitting? If yes, why?</mark>
   High Variance: Unregularized decision trees have high variance, meaning they are sensitive to small fluctuations in the training data. As the tree grows deeper, it becomes more complex and can capture even the noise in the data, leading to a highly irregular decision boundary.

   Complexity: Decision trees can grow to be very complex, especially if there are many features or the tree is not pruned. A complex tree can fit the training data very well but may not generalize well to new, unseen data. This complexity allows the model to capture spurious patterns in the training data that do not reflect the underlying relationships in the population.

   Memorization: Unregularized decision trees have the capacity to memorize the training data, effectively learning each training instance and its associated label. This memorization of the training data can lead to poor generalization performance on new data since the model fails to capture the underlying patterns that are true for the entire population.

6. <mark>What is an ensemble technique in machine learning?</mark>
   Bagging (Bootstrap Aggregating): Bagging involves training multiple instances of the same model on different subsets of the training data, typically using bootstrapping (sampling with replacement). The final prediction is obtained by averaging or taking a majority vote of the predictions made by individual models. Random Forest is a popular ensemble method based on bagging, where decision trees are trained on bootstrapped samples of the data and aggregated to make predictions.

   Boosting: Boosting is an iterative ensemble technique that sequentially trains a series of weak learners (models that perform slightly better than random chance) and assigns higher weights to data points that were misclassified by previous models. The final prediction is a weighted sum of the predictions made by each weak learner. AdaBoost and Gradient Boosting Machines (GBM) are examples of boosting algorithms.

   Stacking (Stacked Generalization): Stacking combines the predictions of multiple base models by training a meta-model on the outputs of these models. Instead of simply averaging or voting on predictions, stacking learns how to best combine the predictions of base models to produce the final prediction. Stacking is typically used with diverse base models to capture different aspects of the data.

   Voting: Voting, also known as majority voting or plurality voting, combines the predictions of multiple models by taking a majority vote or averaging the predicted probabilities across models. It can be used with any combination of models, including models trained on different algorithms or subsets of features.

7. <mark>What is out-of-bag error in random forests?</mark>
   In Random Forests, the out-of-bag (OOB) error is an estimate of the model's performance on unseen data, calculated using the data points that were not included in the bootstrap sample used to train each individual decision tree in the forest.

   Here's how the out-of-bag error estimation works in Random Forests:

Bootstrap Sampling: Each decision tree in the Random Forest is trained on a bootstrap sample of the original training data. This means that for each tree, some data points from the original dataset are sampled with replacement to create a new dataset of the same size. Some data points may be included multiple times in the bootstrap sample, while others may not be included at all.

Out-of-Bag Data: Approximately one-third of the original data points are left out from each bootstrap sample. These left-out data points constitute the out-of-bag (OOB) data for each tree.

Error Estimation: After training each decision tree, the model's performance is evaluated using the out-of-bag data points that were not included in the training sample for that particular tree. These out-of-bag data points serve as a validation set for estimating the model's prediction error.

Aggregation: The out-of-bag error estimates from all the decision trees in the Random Forest are aggregated to obtain the overall out-of-bag error for the entire ensemble model. This aggregated out-of-bag error provides an unbiased estimate of the model's performance on unseen data and can be used to assess the model's generalization ability.

8. What is K-fold cross-validation
   K-fold cross-validation is a resampling technique used to evaluate the performance of a machine learning model. It is particularly useful when the dataset is limited or when there's a need to estimate how well the model will perform on unseen data.

   Here's how K-fold cross-validation works:

   Partitioning the Data: The original dataset is divided into K subsets of approximately equal size. These subsets are often called "folds".

   Training and Testing: The model is trained K times, with each iteration using a different fold as the test set and the remaining K-1 folds as the training set. For example, in the first iteration, the first fold is used as the test set and the remaining K-1 folds are used for training. In the second iteration, the second fold is used as the test set, and so on.

   Evaluation: After each iteration, the performance metric (such as accuracy, precision, recall, or F1-score) is computed using the test set. This process results in K performance estimates, one for each fold.

   Aggregation: The K performance estimates are then averaged to obtain a single performance metric, which represents the overall performance of the model.

   The value of K is typically chosen based on the size of the dataset and computational resources. Common choices for K include 5-fold and 10-fold cross-validation, although other values such as 3-fold and 20-fold can also be used.

   K-fold cross-validation provides several benefits:

It provides a more reliable estimate of the model's performance compared to a single train-test split, as it uses multiple train-test splits.
It allows the entire dataset to be used for both training and testing, maximizing the use of available data.
It helps in identifying potential issues such as overfitting or underfitting by assessing the model's performance across multiple train-test splits.

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

Hyperparameter tuning, also known as hyperparameter optimization, is the process of selecting the optimal set of hyperparameters for a machine learning model. Hyperparameters are parameters that are set prior to training and control the behavior of the learning algorithm. Unlike model parameters, which are learned during training, hyperparameters are not optimized directly from the training data but are instead chosen based on prior knowledge, heuristics, or through a search process.

Hyperparameter tuning is done for several reasons:

Improving Model Performance: Selecting the right hyperparameters can significantly improve the performance of a machine learning model. Hyperparameters influence the complexity, flexibility, and generalization ability of the model. By tuning hyperparameters, we can find the settings that lead to the best performance on a given dataset.

Avoiding Overfitting or Underfitting: Hyperparameters play a crucial role in controlling the complexity of the model. By tuning hyperparameters, we can prevent overfitting (where the model learns to fit the training data too closely and performs poorly on unseen data) or underfitting (where the model is too simple to capture the underlying patterns in the data).

Adapting to Different Datasets: The optimal hyperparameters for a model may vary depending on the characteristics of the dataset, such as its size, dimensionality, and noise level. Hyperparameter tuning allows us to adapt the model to different datasets and maximize its performance across various tasks.

Optimizing Computational Resources: Hyperparameter tuning can help optimize computational resources by finding the most efficient settings for training the model. For example, tuning hyperparameters such as learning rate or batch size can affect the convergence speed and memory usage during training.

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

Divergence: A large learning rate can cause the optimization process to diverge, meaning that instead of converging to the optimal solution, the parameter updates become increasingly large, and the algorithm fails to find a solution. This can happen because the updates are so large that they overshoot the minimum of the loss function, causing the algorithm to oscillate or move away from the optimal solution.

Overshooting the Minimum: With a large learning rate, the gradient descent updates can overshoot the minimum of the loss function, leading to oscillations around the optimal

solution or instability in the optimization process. This results in slower convergence or failure to converge altogether.

Instability: Large learning rates can make the optimization process unstable, especially in regions of the parameter space where the gradient is steep. This can lead to large fluctuations in the parameter values and make it difficult for the algorithm to settle on a stable solution.

Numeric Overflow: In some cases, using a large learning rate can lead to numeric overflow or underflow issues, especially when working with very large or very small parameter values. This can cause the algorithm to produce NaNs (Not-a-Number) or Infinity values, resulting in an unstable optimization process.

Poor Generalization: Large learning rates can lead to poor generalization performance, as the algorithm may fail to find the optimal solution or converge to a suboptimal solution that performs poorly on unseen data. This can result in models that are overly sensitive to small variations in the training data and perform poorly on new data.

==12. Can we use Logistic Regression for classification of Non-Linear Data? If not, why==?
Logistic Regression is a linear classification algorithm, meaning it models the relationship between the input features and the output label using a linear function. Therefore, by its nature, Logistic Regression is limited in its ability to handle non-linear relationships between the features and the target variable.However, there are ways to extend Logistic Regression to handle non-linear data:

Feature Engineering: One approach is to engineer non-linear features from the original features, such as polynomial features or interaction terms. By transforming the input features, you can introduce non-linear relationships into the model and potentially improve its performance.

Kernel Tricks: Another approach is to use kernelized versions of Logistic Regression, such as Kernel Logistic Regression or Support Vector Machines (SVM) with a non-linear kernel (e.g., polynomial kernel, radial basis function kernel). These methods map the input features into a higher-dimensional space where the data may be linearly separable, even if it was not in the original feature space.

Ensemble Methods: Ensemble methods such as Random Forests or Gradient Boosting Machines (GBM) can handle non-linear relationships more effectively than Logistic Regression. By combining multiple weak learners (e.g., decision trees), these methods can capture complex non-linear patterns in the data and provide better classification performance.

Neural Networks: Neural networks, especially deep learning models, are highly flexible and capable of learning complex non-linear relationships in the data. They can be used as a more powerful alternative to Logistic Regression for handling non-linear classification tasks.

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

The bias-variance tradeoff is a fundamental concept in machine learning that relates to the balance between bias and variance in the predictive performance of a model. It describes the tradeoff between the model's ability to capture the underlying patterns in the data (bias) and its sensitivity to small fluctuations or noise in the training data (variance).

Bias: Bias refers to the error introduced by the model's assumptions or simplifications when approximating the true underlying relationship between the features and the target variable. A high bias model tends to underfit the data, meaning it fails to capture the complexity of the underlying patterns and performs poorly on both the training and test data.

Variance: Variance refers to the model's sensitivity to small fluctuations or noise in the training data. A high variance model tends to overfit the data, meaning it captures not only the underlying patterns but also the noise in the training data. As a result, it performs well on the training data but poorly on the test data.

The bias-variance tradeoff arises because reducing bias typically increases variance, and vice versa. Here's how it works:

High Bias, Low Variance: Models with high bias and low variance are simple and have strong assumptions about the data. They tend to underfit the training data but may generalize well to unseen data if the underlying assumptions hold true. Examples include linear models and decision trees with limited depth.

Low Bias, High Variance: Models with low bias and high variance are complex and flexible, often with fewer assumptions about the data. They can capture complex relationships in the training data but are sensitive to noise and fluctuations. They tend to overfit the training data and may not generalize well to unseen data. Examples include deep neural networks and decision trees with large depth.

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

Sure, here's a brief description of Linear, Radial Basis Function (RBF), and Polynomial kernels used in Support Vector Machines (SVM):

Linear Kernel:

The Linear kernel is the simplest kernel used in SVM.
It represents a linear decision boundary between classes.
The decision boundary is a hyperplane that separates the classes by maximizing the margin between them.
It works well when the data is linearly separable.
Radial Basis Function (RBF) Kernel:

The RBF kernel is a popular non-linear kernel used in SVM.
It is effective in capturing complex non-linear relationships between features and target variable.

The RBF kernel transforms the input features into a higher-dimensional space using a Gaussian function.

It allows SVM to model more complex decision boundaries that are not necessarily linear.

The RBF kernel has two hyperparameters: gamma (γ), which controls the width of the Gaussian function, and C, which controls the trade-off between maximizing the margin and minimizing the classification error.

Polynomial Kernel:

The Polynomial kernel is another non-linear kernel used in SVM.

It computes the dot product of the input features in a higher-dimensional space using a polynomial function.

The degree of the polynomial function determines the complexity of the decision boundary.

Like the RBF kernel, the Polynomial kernel allows SVM to capture non-linear relationships between features and target variable.

It has a hyperparameter, degree, which controls the degree of the polynomial function.

The Polynomial kernel can be less computationally expensive compared to the RBF kernel for some datasets.