# MACHINE LEARNING WORK SHEET SET 05

***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 (R2) and Residual Sum of Squares (RSS) are both measures of the goodness of fit of a regression model, but they capture different aspects of the model's performance.

R-squared measures the proportion of variance in the dependent variable (i.e., the outcome variable) that is explained by the independent variables (i.e., the predictors) in the model. R2 ranges from 0 to 1, with higher values indicating better fit. A high R2 suggests that the model is able to explain a large amount of the variability in the dependent variable.

On the other hand, RSS measures the total amount of unexplained variance in the dependent variable. RSS is calculated as the sum of the squared residuals (i.e., the differences between the predicted values and the actual values) and reflects the amount of variability in the dependent variable that is not accounted for by the independent variables.

In general, a good model should have both a high R2 and a low RSS. However, if we have to choose between the two, R2 is often considered a better measure of goodness of fit than RSS, because it provides a measure of the proportion of variance in the dependent variable that is explained by the model. R2 is also more easily interpretable than RSS, as it is on a scale of 0 to 1.

However, it's worth noting that R2 can be misleading if used alone. For example, a high R2 does not necessarily indicate a good model if the predictors are not meaningful or relevant to the dependent variable. Similarly, a low R2 does not necessarily mean a bad model if the predictors are meaningful and relevant, but the dependent variable is highly variable and difficult to predict. In these cases, RSS may be a better measure to consider.

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, Total Sum of Squares (TSS), Explained Sum of Squares (ESS), and Residual Sum of Squares (RSS) are all measures of the variability in the dependent variable (i.e., the outcome variable) that can be attributed to different sources.

TSS is the total amount of variability in the dependent variable, and can be decomposed into the variability that is explained by the independent variables (i.e., ESS) and the variability that is not explained by the independent variables (i.e., RSS).

ESS is the variability in the dependent variable that is explained by the independent variables in the model. It represents the sum of the squared differences between the predicted values and the mean of the dependent variable. In other words, it measures how much of the total variability in the dependent variable is explained by the model.

RSS is the variability in the dependent variable that is not explained by the independent variables in the model. It represents the sum of the squared differences between the actual values and the predicted values. In other words, it measures how much of the total variability in the dependent variable is left unexplained by the model.

The equation relating TSS, ESS, and RSS is:

$$TSS = ESS + RSS$$

This equation reflects the fact that the total variability in the dependent variable can be decomposed into the variability that is explained by the model and the variability that is not explained by the model. In other words, the sum of the variability that is explained and the variability that is not explained equals the total variability.

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

**Answer:**

Regularization is a technique used in machine learning to prevent overfitting of the model. Overfitting occurs when the model learns the noise in the training data, resulting in poor generalization performance on new, unseen data.

Regularization methods add a penalty term to the objective function of the learning algorithm to discourage the model from learning complex or intricate patterns in the training data that may not generalize well. The penalty term is designed to reduce the magnitude of the weights or coefficients of the model, thereby simplifying the model and reducing its capacity to fit the noise in the training data.

The most common regularization methods used in machine learning are L1 regularization, L2 regularization, and dropout regularization. L1 regularization adds a penalty proportional to the absolute value of the weights, while L2 regularization adds a penalty proportional to the square of the weights. Dropout regularization randomly drops out some neurons during training to prevent them from co-adapting too much to the input.

Overall, regularization is necessary in machine learning to improve the generalization performance of the model by preventing overfitting and improving its ability to make accurate predictions on new, unseen data.

### 4. What is Gini–impurity index?

**Answer:**

The Gini impurity index is a measure of impurity used in decision tree algorithms for classification tasks. It is a statistical measure that quantifies the degree of probability of incorrectly classifying a randomly chosen element in a dataset.

The Gini impurity index measures the probability of misclassification of a randomly chosen element if it is randomly assigned to a class according to the class distribution in the dataset. The Gini impurity index is defined as follows:

$$\text{Gini Index} = 1 - \Sigma \, (p\_i)^2$$

where $p\_i$ is the probability of an element belonging to a particular class i.

The Gini index takes on a value between 0 and 1, with 0 indicating a perfect classification, and 1 indicating complete uncertainty in the classification. The lower the Gini index, the better the separation of the data into different classes. Therefore, decision trees aim to minimize the Gini index when splitting a node by selecting the feature that maximally reduces the Gini index.

In summary, the Gini impurity index is a measure of the impurity of a set of elements that quantifies the probability of misclassification of a randomly chosen element. It is used in decision tree algorithms to select the best feature to split a node and obtain a better classification.

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

**Answer:**

Yes, unregularized decision trees are prone to overfitting because they can learn the noise in the training data, resulting in high variance and poor generalization performance on new, unseen data.

Decision trees are a non-parametric model that can learn complex decision boundaries by recursively partitioning the feature space. Each split in the decision tree selects the best feature and threshold that maximally reduces the impurity of the data at that node. As a result, decision trees can easily overfit the training data by fitting the noise instead of the underlying patterns.

Unregularized decision trees continue to grow until they have perfectly classified the training data, resulting in a large, deep tree with many leaves. This can lead to overfitting, where the model becomes too complex and specific to the training data, making it difficult to generalize to new, unseen data.

To prevent overfitting, regularization techniques can be used, such as pruning the tree, setting a minimum number of samples required to split a node, or limiting the maximum depth of the tree. These regularization techniques can simplify the decision tree and reduce its capacity to fit the noise in the data, resulting in better generalization performance.

**6.  What is an ensemble technique in machine learning?**

**Answer:**

Ensemble techniques in machine learning are a type of meta-algorithm that combine multiple individual models to improve the accuracy and robustness of the predictions. Ensemble methods are based on the idea that a group of models will perform better than any single model, as each model may capture different aspects of the data and compensate for each other's weaknesses.

There are several types of ensemble methods, including:

1.  Bagging: Bootstrap Aggregation or Bagging is a technique that trains multiple models on different subsets of the training data with replacement. The predictions from all the models are combined using averaging or voting to make a final prediction.

2.  Boosting: Boosting is a technique that trains multiple models sequentially, where each subsequent model tries to improve the performance of the previous model by focusing on the misclassified samples. The predictions from all the models are combined using weighted averaging to make a final prediction.

3. Stacking: Stacking is a technique that trains multiple models of different types, such as decision trees, neural networks, or support vector machines, and uses the output of each model as input to a higher-level model, which then makes the final prediction.

4. Random Forest: A random forest is an ensemble of decision trees that uses bagging and random feature selection to improve the performance of the individual trees. Each tree is trained on a different subset of the data with replacement, and at each node, a random subset of the features is considered for splitting.

Ensemble methods are widely used in machine learning because they can improve the accuracy and robustness of the predictions, and they are less prone to overfitting than single models. Ensemble methods can be applied to both classification and regression tasks, and they are particularly useful when the individual models have high variance or bias.

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

*Answer:*

The main difference between Bagging and Boosting techniques in machine learning is in how they create the ensemble of models from the training data and how they combine the predictions of the individual models.

Bagging (Bootstrap Aggregation) and Boosting are both ensemble methods that use multiple models to improve the accuracy and robustness of the predictions. However, they differ in the following ways:

1. **Sample selection:** Bagging selects random subsets of the training data with replacement to train each model independently. In contrast, Boosting selects samples according to their misclassification rate to focus on the difficult samples and improve the model's performance.

2. **Model weights:** Bagging uses equal weights for each model in the ensemble, whereas Boosting assigns higher weights to the better-performing models and lower weights to the poorer-performing models.

3. **Prediction aggregation:** In Bagging, the predictions of each model are aggregated by taking the average or voting across all the models. In Boosting, the predictions of each model are weighted by the model's performance and combined using a weighted average.

4. **Final prediction:** In Bagging, the final prediction is made by taking the average or voting across all the models. In Boosting, the final prediction is made by summing the weighted predictions of each model.

In summary, Bagging and Boosting are both powerful ensemble methods that can improve the accuracy and robustness of the predictions. Bagging focuses on reducing the variance of the model by averaging over multiple models, while Boosting focuses on reducing the bias of the model by emphasizing the difficult samples.

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

*Answer:*

Out-of-bag (OOB) error in random forests is an estimate of the generalization error of the model. The OOB error is calculated by using the samples that were not selected in the bootstrap sample to evaluate the performance of the model.

In a random forest, each decision tree is trained on a random subset of the training data with replacement, and the remaining samples that are not selected in the bootstrap sample are called OOB samples. These OOB samples are not used in the training of the decision tree and can be used to estimate the performance of the model.

To calculate the OOB error, each OOB sample is passed through all the trees in the forest, and the class with the highest frequency is assigned to the sample. The predicted class is then compared to the true class label of the sample, and the error rate is calculated as the proportion of misclassified samples.

The OOB error provides a useful estimate of the generalization error of the random forest model, as it uses the samples that were not used in the training of the individual trees. The OOB error can also be used to tune the hyperparameters of the random forest, such as the number of trees in the forest, the maximum depth of the trees, and the number of features considered at each split.

Overall, the OOB error in random forests is a powerful tool for evaluating the performance of the model and tuning the hyperparameters for optimal performance.

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

*Answer:*

K-fold cross-validation is a technique used to evaluate the performance of a machine learning model by partitioning the available data into K subsets or "folds". The basic idea behind k-fold cross-validation is to divide the data into K equally sized subsets, and then use K-1 subsets as the training data and the remaining subset as the validation data. This process is repeated K times, with each subset being used exactly once as the validation data.

The general steps involved in performing k-fold cross-validation are:

1. Shuffle the data randomly.

2. Split the data into K equally sized folds.

3. For each fold, use that fold as the validation data and the remaining folds as the training data.

4. Train the model on the training data and evaluate it on the validation data.

5. Record the performance metric (such as accuracy, F1-score, etc.) of the model.

6. Repeat steps 3-5 for each fold.

7. Calculate the average performance metric across all the folds.

The advantage of k-fold cross-validation is that it allows us to use all of the available data for both training and validation, while also reducing the variance of the performance estimate. It also helps to reduce overfitting by evaluating the model on multiple independent subsets of the data.

The value of K is typically chosen to be 5 or 10, but it can also be adjusted based on the size of the dataset and the computational resources available. In practice, k-fold cross-validation is often used to tune the hyperparameters of a model, such as the regularization strength or the learning rate, to obtain the best performance.

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

*Answer:*

Hyperparameter tuning in machine learning refers to the process of selecting the optimal values for the hyperparameters of a machine learning algorithm. Hyperparameters are parameters that are not learned from the data during training, but rather are set before training and control the behavior of the algorithm.

The hyperparameters of a machine learning algorithm can have a significant impact on its performance, and finding the optimal values for these hyperparameters is often a crucial step in building an accurate and robust model.

Hyperparameter tuning is typically done using a validation set or by using cross-validation. The general process involves selecting a range of values for each hyperparameter, training the model on the training set using different combinations of hyperparameter values, and evaluating the performance of the model on a validation set or using cross-validation. The optimal combination of hyperparameters is then selected based on the performance metric of interest, such as accuracy, F1-score, or mean squared error.

Some common hyperparameters that are often tuned in machine learning algorithms include the learning rate, regularization strength, number of hidden units in a neural network, kernel parameters in a support vector machine, number of trees in a random forest, and depth of a decision tree.

Hyperparameter tuning is done to optimize the performance of a machine learning algorithm and prevent overfitting. By selecting the optimal values for the hyperparameters, the algorithm can generalize better to new data and produce more accurate predictions.

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

*Answer:*

If the learning rate in gradient descent is too large, it can lead to the following issues:

1. Divergence: A large learning rate can cause the algorithm to diverge, meaning that the updates to the weights become larger and larger and the optimization process becomes unstable. This can result in the loss function oscillating or even increasing over time, leading to poor model performance.

2. Overshooting: A large learning rate can cause the algorithm to overshoot the optimal weights and bounce back and forth across the optimal point, resulting in a slower convergence rate and longer training time.

3. Local Minima: A large learning rate can cause the algorithm to jump over the global minimum and get stuck in a local minimum. This is because a large learning rate can cause the algorithm to skip over the global minimum and converge to a suboptimal solution.

4. Unstable convergence: A large learning rate can cause the optimization process to be unstable and oscillate around the optimal solution, leading to slow convergence or failure to converge altogether.

To avoid these issues, it is important to choose an appropriate learning rate that balances the trade-off between convergence speed and stability. In practice, a common approach is to use a

learning rate schedule that decreases the learning rate over time, such as a learning rate decay or adaptive learning rate methods like Adam or RMSprop, to ensure stable convergence and optimal performance.

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

**Answer:**

Logistic regression is a linear classification algorithm that models the relationship between the input features and the binary output variable using a linear function. Therefore, it is not suitable for classification of non-linear data that cannot be separated by a linear decision boundary.

If the data is non-linear, using logistic regression may lead to underfitting, where the model is not able to capture the complexity of the data and has a high bias. In such cases, using a linear decision boundary may not be sufficient to accurately classify the data.

To classify non-linear data, we can use non-linear classification algorithms such as decision trees, random forests, support vector machines (SVMs), or neural networks. These algorithms are capable of learning non-linear decision boundaries and can accurately classify non-linear data.

In summary, while logistic regression is a powerful and widely used classification algorithm, it is not suitable for non-linear data. For non-linear data, we need to use non-linear classification algorithms that can model complex decision boundaries.

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

**Answer:**

Adaboost (short for Adaptive Boosting) and Gradient Boosting are both ensemble learning techniques that combine multiple weak learners to form a strong learner. However, they differ in their approach to combining these weak learners and updating the weights of training examples.

Adaboost:

- Adaboost combines a set of weak learners (typically decision trees with a single split, also known as decision stumps) and assigns weights to each example in the training set.

- The algorithm iteratively trains the weak learners on a modified version of the training data, where the weights of misclassified examples are increased and the weights of correctly classified examples are decreased.

- In each iteration, the weak learner is trained to minimize the weighted error of the previous iteration, with more weight given to examples that were misclassified.

- The final prediction is a weighted sum of the predictions of all the weak learners.

Gradient Boosting:

- Gradient Boosting also combines a set of weak learners, but they are typically more complex than decision stumps, such as decision trees with multiple splits.

- The algorithm iteratively trains the weak learners on the residuals (the difference between the predicted and actual values) of the previous iteration.

- In each iteration, the weak learner is trained to predict the residuals, and the weights of the training examples are not modified.

- The final prediction is the sum of the predictions of all the weak learners.

In summary, while Adaboost and Gradient Boosting are both ensemble learning techniques that use multiple weak learners to form a strong learner, they differ in how they combine these weak learners and update the weights of training examples. Adaboost assigns weights to examples and minimizes the weighted error of the previous iteration, while Gradient Boosting trains the weak learners on the residuals of the previous iteration.

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

*Answer:*

Bias-variance tradeoff is a fundamental concept in machine learning that describes the relationship between the complexity of a model and its ability to generalize to new, unseen data. It refers to the tradeoff between the bias and variance of a model.

Bias measures how much the predictions of a model differ from the true values, on average, when trained on different data sets. A model with high bias makes strong assumptions about the data and tends to underfit, or oversimplify, the relationships between the input and output variables.

Variance measures how much the predictions of a model differ from the true values, on average, when trained on different data sets. A model with high variance is sensitive to noise and tends to overfit, or memorize, the training data.

The goal of machine learning is to find a model with low bias and low variance that can generalize well to new data. However, decreasing one often results in an increase in the other, creating a tradeoff between the two.

For example, a simple model with few parameters may have high bias and low variance, while a complex model with many parameters may have low bias and high variance. The key is to find the right balance between the two, by choosing a model with an appropriate level of complexity and using regularization techniques such as L1 or L2 regularization, dropout, or early stopping to reduce the variance of the model.

In summary, the bias-variance tradeoff describes the balance between the ability of a model to capture the complexity of the data and its ability to generalize to new, unseen data. The goal is to find a model with an appropriate level of complexity that minimizes both bias and variance.

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

*Answer:*

SVM (Support Vector Machine) is a popular supervised learning algorithm used for classification and regression analysis. SVMs work by transforming the input data into a higher dimensional space, where it can be separated into classes by a hyperplane. Kernels are used in SVMs to compute the inner product between pairs of data points in the transformed space. There are different types of kernels available, including:

1. **Linear kernel:** The linear kernel is the simplest kernel used in SVM. It computes the inner product between pairs of data points in the original space, which is equivalent to projecting

the data onto a higher dimensional space using a linear transformation. It works well when the data is linearly separable, meaning it can be separated by a straight line or hyperplane.

2. **RBF (Radial Basis Function) kernel:** The RBF kernel is a popular choice for SVM. It maps the data into an infinite dimensional space using a Gaussian function, which allows it to capture complex, non-linear relationships between the input and output variables. The RBF kernel is suitable for data that is not linearly separable, and it can also handle outliers in the data.

3. **Polynomial kernel:** The polynomial kernel is another non-linear kernel used in SVM. It maps the data into a higher dimensional space using a polynomial function, which allows it to capture non-linear relationships between the input and output variables. The polynomial kernel is suitable for data that has polynomial relationships between the input and output variables, and it can also handle outliers in the data.

In summary, SVM is a powerful machine learning algorithm that can be used for classification and regression analysis. Kernels are used in SVM to transform the input data into a higher dimensional space, where it can be separated into classes by a hyperplane. The linear kernel is the simplest kernel used in SVM, while the RBF and polynomial kernels are used for non-linear data.