# Assignment- 5

## Statistics

**Batch DS2311**

Name: Alvin V Regi

1. **Using a goodness of fit,we can assess whether a set of obtained frequencies differ from a set of frequencies.**

**Answer-**

Expected

A goodness of fit test is a statistical method that compares the observed frequencies of a categorical variable with the expected frequencies based on a hypothesized distribution. It can be used to determine whether the observed data are consistent with the assumed distribution or not.

2. **Chisquare is used to analyse**

**Answer-**

Frequencies.

Chi-square is a statistical test that is used to analyse the relationship between categorical variables. It can be used to test whether the observed frequencies of different categories are equal to the expected frequencies, or whether the frequencies of different categories vary across groups.

3. **What is the mean of a Chi Square distribution with 6 degrees of freedom?**

**Answer-**

6.

The mean of a chi-square distribution with k degrees of freedom is equal to k.

4. **Which of these distributions is used for a goodness of fit testing?**

**Answer-**

Chi-squared distribution.

A goodness of fit test compares the observed frequencies of different categories with the expected frequencies based on a hypothesized distribution. The chi-square test is one of the most common goodness of fit tests, and it uses the chi-square distribution to calculate the p-value of the test.

5. **Which of the following distributions is Continuous?**

**Answer-**

F Distribution.

A continuous distribution is a probability distribution that can take any value within a specified range or interval. The F distribution is a continuous distribution that is used to compare the variances of two populations or to test the significance of regression models.

6. **A statement made about a population for testing purpose is called?**

**Answer-**

Hypothesis.

A hypothesis is a statement made about a population parameter, such as the mean or the proportion, that can be tested using sample data. A hypothesis test is a statistical procedure that evaluates the evidence in favor or against the hypothesis and provides a conclusion based on the level of significance and the test statistic.

7. **If the assumed hypothesis is tested for rejection considering it to be true is called?**

**Answer-**

Null Hypothesis.

A null hypothesis is the hypothesis that is assumed to be true and tested for rejection using sample data. It usually states that there is no difference or no relationship between the population parameters of interest.

8. **If the Critical region is evenly distributed then the test is referred as?**

**Answer-**

Two tailed.

A two-tailed test is a hypothesis test that has two critical regions, one in each tail of the sampling distribution. It is used when the alternative hypothesis is not directional, meaning that it does not specify whether the population parameter is greater than or less than the null value.

9. **Alternative Hypothesis is also called as?**

**Answer-**

Research Hypothesis.

A research hypothesis is the hypothesis that is proposed by the researcher as an alternative to the null hypothesis. It usually states that there is a difference or a relationship between the population parameters of interest.

10. **In a Binomial Distribution, if 'n' is the number of trials and 'p' is the probability of success, then the mean value is given by ___**

**Answer-**

np.

The mean value of a binomial distribution with n trials and p probability of success is given by the formula $\mu = np$.

# Machine Learning

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-**

The R-squared (coefficient of determination) and Residual Sum of Squares (RSS) are both measures of goodness of fit in regression models, but they capture different aspects.

R-squared represents the proportion of the dependent variable's variance that is explained by the independent variables. It ranges from 0 to 1, with 1 indicating a perfect fit. R-squared is advantageous because it provides an intuitive understanding of how well the independent variables explain the dependent variable.

On the other hand, RSS measures the overall difference between the observed and predicted values of the dependent variable. It quantifies the total error or residual in the model. RSS is advantageous because it evaluates the absolute fit of the model and can be used to compare different models.

In terms of which measure is better, it depends on the context and the goals of the analysis. R-squared is commonly used because it provides a clear interpretation of the explanatory power of the model. However, it can be misleading if the model is overfit or if it includes irrelevant variables. In such cases, RSS can be a better measure because it directly assesses the model's predictive accuracy and accounts for the number of variables used.

In summary, R-squared is useful to understand the explanatory power of the model, while RSS is useful to assess the overall fit and compare different models. Both measures have their merits, and researchers often consider them together to gain a comprehensive understanding of the model's goodness of fit.

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 analysis, TSS (Total Sum of Squares), ESS (Explained Sum of Squares), and RSS (Residual Sum of Squares) are important metrics that help in understanding the goodness of fit of the model.

TSS represents the total variation in the dependent variable (y) and is calculated by summing the squared differences between each observed y value and the mean of y.

ESS represents the variation in the dependent variable (y) that is explained by the independent variables (x) in the regression model. It can be calculated by summing the squared differences between the predicted y values (ŷ) and the mean of y.

RSS represents the remaining variation in the dependent variable (y) that is not explained by the independent variables (x) in the model. It is calculated by summing the squared differences between the observed y values and the predicted y values.

The relationship between these three metrics can be expressed by the equation:

TSS = ESS + RSS

This equation highlights that the total sum of squares (TSS) can be decomposed into the explained sum of squares (ESS) and the residual sum of squares (RSS). The ESS and RSS together make up the total variation in the dependent variable. The higher the ESS and the lower the RSS, the better the model's goodness of fit.

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

**Answer-**

Regularization is a technique used in machine learning to prevent overfitting and improve the generalization ability of models. Overfitting occurs when a model performs well on the training data but fails to generalize to unseen data.

Regularization addresses this issue by adding a penalty term to the loss function during model training. This penalty discourages complex models that may fit the training data perfectly but have poor performance on new data. Regularization helps in achieving a balance between model complexity and generalization.

The need for regularization arises from the following reasons:

1. Preventing overfitting: Regularization techniques such as L1 and L2 regularization constrain the model's parameters, preventing them from becoming too large or too complex. This helps in reducing overfitting by discouraging the model from memorizing noise or irrelevant patterns in the training data.

2. Handling multicollinearity: In cases where the input features are highly correlated (multicollinearity), regularization techniques help in reducing the impact of such correlations on the model. This improves the stability and interpretability of the model.

3. Feature selection: Regularization can drive certain model coefficients to zero, effectively performing feature selection. This is particularly useful when dealing with high-dimensional datasets, where it helps in identifying and focusing on the most relevant features.

4. Improving model interpretability: Regularization can simplify the model by reducing the number of features or shrinking their coefficients. This leads to a more interpretable model, as it focuses on the most important features and reduces the impact of noise.

Overall, regularization techniques help in improving the performance, generalization, stability, and interpretability of machine learning models, making them more robust and reliable in real-world scenarios.

### 4. What is Gini–impurity index?

**Answer-**

The Gini impurity index is a measure used in decision tree algorithms to evaluate the impurity or heterogeneity of a set of class labels. It quantifies the probability of incorrectly classifying a randomly chosen element in a dataset based on the distribution of class labels.

A Gini impurity index of 0 indicates a pure node where all elements belong to the same class. A Gini impurity index of 1 indicates maximum impurity where the elements are evenly distributed among different classes.

In the context of decision trees, the Gini impurity index is used as a criterion to determine the optimal split at each node. The split that minimizes the weighted sum of the Gini impurity indices for the resulting child nodes is chosen as the best split.

By repeatedly splitting the data based on the Gini impurity index, decision trees can effectively partition the data into homogeneous subsets, ultimately leading to accurate and reliable predictions.

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

**Answer-**

Yes, unregularized decision trees are prone to overfitting. There are several reasons why this is the case:

1. Capturing noise: Decision trees have the ability to learn intricate patterns and details from the training data, including noise and outliers. Without any form of regularization, the tree can become overly complex and fit the noise in the data, leading to poor generalization on unseen data.

2. Memorizing the training data: Unregularized decision trees have the potential to memorize the training data by creating branches and leaves to perfectly fit each training example. This results in a high degree of complexity and specificity, which may not generalize well to new, unseen data.

3. Lack of constraint on tree depth: Unregularized trees have no restrictions on their depth or the number of splits they can make. This can lead to trees that are too deep, with numerous decision nodes and branches, making them highly specific to the training data and less likely to generalize to new data.

4. Sensitivity to small changes in data: Unregularized decision trees can be sensitive to small changes in the training data, resulting in different splits and structures in the tree. This instability can lead to overfitting, as the tree may fit noise or irrelevant patterns that are specific to the particular training set.

To mitigate overfitting in decision trees, various regularization techniques can be applied, such as pruning, limiting tree depth, setting minimum sample requirements for splitting, or using ensemble methods like random forests or gradient boosting. These techniques help control the complexity of the tree and improve its generalization capabilities.

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

**Answer-**

Ensemble techniques in machine learning involve combining multiple individual models to create a more accurate and robust predictive model. The idea behind ensemble methods is that by aggregating the predictions of multiple models, the strengths of each model can compensate for the weaknesses of others, leading to improved overall performance.

There are two main types of ensemble techniques:

1. Bagging (Bootstrap Aggregating): In bagging, multiple models, such as decision trees, are trained independently on different subsets of the training data, created through bootstrap sampling. The predictions of these models are then combined through majority voting (classification) or averaging (regression) to make the final prediction.

2. Boosting: Boosting is an iterative ensemble technique where models are trained sequentially, with each subsequent model focusing on the instances that were misclassified by the previous models. Boosting algorithms, like AdaBoost and Gradient Boosting, assign higher weights to the misclassified instances and adjust the model's parameters to improve their predictions.

Ensemble techniques offer several benefits, including improved accuracy, better generalization ability, reduced overfitting, and increased stability. They are widely used in various machine learning tasks, such as classification, regression, and anomaly detection, and have proven to be highly effective in improving model performance.

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

**Answer-**

The main differences between bagging and boosting techniques are as follows:

1. Training process: In bagging, each model in the ensemble is trained independently on different subsets of the training data through bootstrap sampling. The models are unrelated and trained in parallel. In boosting, the models are trained sequentially, with each subsequent model focusing on the instances that were misclassified by the previous models. The models are trained iteratively, and the training process is influenced by the performance of the previous models.

2. Weighting of instances: In bagging, each model is trained on a random subset of the training data, and all instances have equal weights. In boosting, instances are assigned weights, with higher weights given to the misclassified instances. This allows subsequent models to focus more on the difficult instances and improve their predictions.

3. Combination of predictions: In bagging, the predictions of the individual models are combined through majority voting (classification) or averaging (regression) to make the final prediction. In boosting, the predictions are combined through weighted voting, where each model's prediction is weighted based on its performance during training.

4. Handling of errors: Bagging reduces the variance of the model by averaging multiple predictions, but it does not explicitly correct or reduce errors made by individual models. Boosting, on the other hand, focuses on minimizing errors by iteratively adjusting the model's parameters and assigning higher weights to misclassified instances.

Overall, bagging aims to reduce variance and improve stability by averaging multiple models, while boosting focuses on reducing bias and improving accuracy by iteratively training models that correct the errors made by previous models.

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

**Answer-**

The out-of-bag (OOB) error is a method used to estimate the performance of a random forest model without the need for cross-validation or a separate validation set. In random forests, each decision tree is trained using a bootstrap sample of the original training data, which means that some instances are left out or "out-of-bag" in each bootstrap sample.

The OOB error is then calculated by evaluating each instance on the decision trees that were not trained using that particular instance. The predictions from these out-of-bag instances are compared to their true labels, and the error is calculated. The OOB error is the average error over all out-of-bag instances and serves as an estimate of the model's generalization error.

The OOB error provides a reliable estimate of how well the random forest model will perform on unseen data and is useful for model evaluation and hyperparameter tuning.

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

**Answer-**

K-fold cross-validation is a technique used to evaluate the performance and generalization ability of a machine learning model. It involves dividing the available labeled data into K subsets or folds.

The process of K-fold cross-validation can be summarized as follows:

1. The data is divided into K equal-sized folds.

2. The model is trained K times, with each iteration using K-1 folds as the training data and the remaining fold as the validation data.

3. The performance of the model is evaluated on each validation set, resulting in K performance metrics.

4. The final performance metric is usually calculated as the average or the aggregated result of the K performance metrics.

The main advantage of K-fold cross-validation is that it allows for a more robust estimation of the model's performance by using different subsets of data for training and validation. It helps to mitigate the impact of data variability and provides a more accurate evaluation of the model's ability to generalize to unseen data.

Common variations of K-fold cross-validation include stratified K-fold, where the class distribution is preserved in each fold, and repeated K-fold, where the process is repeated multiple times with different random splits.

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

**Answer-**

Hyperparameter tuning involves the process of selecting the optimal values for hyperparameters in a machine learning model. Hyperparameters are parameters that are not learned from data but are set prior to training. Examples of hyperparameters include the learning rate in a neural network, the depth of a decision tree, or the number of neighbors in a k-nearest neighbors algorithm.

Hyperparameter tuning is done to find the combination of hyperparameter values that results in the best performance of the model. The goal is to improve the model's generalization ability, accuracy, and robustness. By systematically tuning hyperparameters, we can find the configuration that minimizes the model's error or maximizes its performance metric, such as accuracy or F1 score.

Hyperparameter tuning can be done using various techniques, including grid search, random search, and Bayesian optimization. It involves evaluating the model's performance on a validation set or using techniques like cross-validation. By finding the optimal hyperparameter values, we can build machine learning models that are more effective and better suited to the specific problem and data at hand.

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

**Answer-**

If we have a large learning rate in Gradient Descent, it can lead to several issues:

1. Overshooting: A large learning rate can cause the algorithm to take large steps in the direction of steepest descent. This can result in overshooting the optimal solution, leading to instability and divergence. The algorithm may fail to converge or oscillate around the minimum instead of converging.

2. Slow convergence: While a large learning rate may initially result in quick progress, it can later cause the algorithm to bounce around the minimum or get stuck in a suboptimal region. This can slow down the convergence process and prevent the algorithm from reaching the global minimum.

3. Unstable gradients: Large learning rates can cause unstable and erratic gradients, especially in deep neural networks. This phenomenon, known as the exploding gradient problem, can make it difficult for the model to learn effectively and can lead to numerical instability during training.

4. Loss function fluctuations: A large learning rate can cause the loss function to fluctuate significantly during training. These fluctuations can make it challenging to determine if the model is making progress or if it is stuck in a local minimum.

To mitigate these issues, it is crucial to choose an appropriate learning rate. Techniques such as learning rate schedules, adaptive learning rate methods (e.g., Adam, RMSprop), and early stopping can be employed to prevent the problems associated with large learning rates and ensure stable and efficient convergence.

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

**Answer-**

Logistic Regression is a linear classification algorithm that assumes a linear relationship between the features and the target variable. It works well for problems with linearly separable data, where a straight line can effectively separate the classes.

However, for non-linear data, Logistic Regression may not be able to capture the underlying patterns and make accurate predictions. Since it is a linear model, it can only learn linear decision boundaries. If the data is non-linearly separable, Logistic Regression may struggle to find a suitable decision boundary and result in poor classification performance.

To handle non-linear data, more complex models such as Support Vector Machines (SVM), Decision Trees, Random Forests, or Neural Networks can be used. These models are capable of learning non-linear relationships and can capture more intricate patterns in the data.

In situations where the non-linear data can be transformed into a linearly separable form, Logistic Regression can still be used by applying non-linear feature engineering techniques such as polynomial features or using kernel tricks to transform the data into a higher-dimensional space where it becomes linearly separable.

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

**Answer-**

Adaboost and Gradient Boosting are both ensemble learning methods that combine multiple weak learners to create a strong predictive model. However, they differ in their approach and the way they update the model during training.

1. Training Process:

   o  Adaboost (Adaptive Boosting): Adaboost starts by assigning equal weights to all training instances. It then sequentially trains weak learners, where each subsequent learner focuses on the instances that were misclassified by the previous learners. The weights of the misclassified instances are increased, allowing subsequent learners to pay more attention to them and improve the overall model performance.

   o  Gradient Boosting: Gradient Boosting, on the other hand, trains weak learners in an additive manner. It starts with an initial prediction, usually the mean or a constant value, and then fits subsequent learners to the residuals (the difference between the actual target values and the predicted values). Each learner is trained to minimize the loss function by taking a step in the direction of the negative gradient of the loss with respect to the predictions of the previous learners. The predictions of all the learners are then added together to make the final prediction.

2. Weak Learners:

   o  Adaboost: Adaboost can work with any weak learner, typically decision stumps (a decision tree with a single split). It focuses on improving the performance of the weak learners by emphasizing the misclassified instances.

   o  Gradient Boosting: Gradient Boosting is not restricted to a specific weak learner and can work with a variety of models, such as decision trees or even linear models. It focuses on creating a sequence of models that gradually reduce the overall error by optimizing the loss function.

3. Handling of Errors:

   o  Adaboost: Adaboost assigns higher weights to misclassified instances, allowing subsequent learners to focus on them and improve their performance. It gives more importance to the difficult instances and tries to minimize the overall training error.

   o  Gradient Boosting: Gradient Boosting learns from the residuals of the previous learners and tries to fit subsequent learners to these residuals. It focuses on reducing the overall loss function and directly optimizes the difference between the predicted and actual values.

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

**Answer-**

The bias-variance trade-off in machine learning refers to the relationship between a model's ability to capture the complexity of the underlying data and its generalization performance.

Bias refers to the error introduced by approximating a real-world problem with a simplified model. A high bias model tends to underfit the data, meaning it oversimplifies the relationships and fails to capture important patterns and details.

Variance, on the other hand, refers to the model's sensitivity to the noise or randomness in the training data. A high variance model tends to overfit the data, meaning it fits the training data too closely and fails to generalize well to unseen data.

The trade-off occurs because reducing bias often increases variance, and vice versa. Finding the right balance is crucial for building a model that generalizes well to unseen data. This balance can be

achieved through techniques like regularization, feature selection, and model selection, or by using ensemble methods like random forests or gradient boosting that combine multiple models to reduce both bias and variance.

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

**Answer-**

1.  Linear Kernel: The linear kernel is the simplest and most straightforward kernel used in SVM. It represents a linear decision boundary, which is a straight line in 2D and a hyperplane in higher dimensions. It is suitable for linearly separable data and performs well when the classes can be separated by a straight line or plane.

2.  RBF (Radial Basis Function) Kernel: The RBF kernel is a popular choice for SVM as it can capture non-linear decision boundaries. It uses a Gaussian function to measure the similarity between data points. It maps the data into a higher-dimensional space where it becomes linearly separable. The RBF kernel is versatile and can handle a wide range of data distributions, making it suitable for various classification tasks.

3.  Polynomial Kernel: The polynomial kernel is used to capture polynomial relationships between data points. It maps the data into a higher-dimensional space using polynomial functions, enabling SVM to learn non-linear decision boundaries in a feature space. The degree of the polynomial determines the complexity of the decision boundary. Lower degrees result in simpler boundaries, while higher degrees can capture more intricate patterns. The polynomial kernel is useful for handling data with non-linear relationships but may be sensitive to the choice of the degree parameter.