ASSIGNMENT - 5

# MACHINE LEARNING

**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 ofgoodness of fit model in regression and why?

**ANS :- The residual sum of squares (RSS) is a statistical technique used to measure the amount of variance in a data set that is not explained by a regression model itself. Instead, it estimates the variance in the residuals, or error term.**

**R-Squared (R² or the coefficient of determination) is a statistical measure in a regression model that determines the proportion of variance in the dependent variable that can be explained by the independent variable. In other words, r-squared shows how well the data fit the regression model (the goodness of fit).24**

2.What are TSS (Total Sum of Squares), ESS (Explained Sum of Squares) and RSS (Residual Sumof Squares) in regression. Also
Also mention the equation relating these three metrics with each other.
**ANS :-**
**1.TSS :- The total sum of squares (TSS) measures how much variation there is in the observed data, while the residual sum of squares measures the variation in the error between the observed data and modeled values.**
**2.ESS :- The explained sum of squares (ESS) is the sum of the squares of the deviations of the predicted values from the mean value of a response variable, in a standard regression model — for example, yi = a + b1x1i + b2x2i + ... + εi, where yi is the i th observation of the response variable, xji is the i th observation of the j th explanatory variable, a and bj are coefficients, i indexes the observations from 1 to n, and εi is the i th value of the error term. In general, the greater the ESS, the better the estimated model performs.**
**3.RSS :- The residual sum of squares (RSS) is a statistical technique used to measure the amount of variance in a data set that is not explained by a regression model itself. Instead, it estimates the variance in the residuals, or error term.**

3.What is the need of regularization in machine learning?
**ANS :- Regularization refers to techniques that are used to calibrate machine learning models in order to minimize the adjusted loss function and prevent overfitting or underfitting. Using Regularization, we can fit our machine learning model appropriately on a given test set and hence reduce the errors in it.**

4.What is Gini–impurity index?

**ANS :- Gini Impurity is a measurement used to build Decision Trees to determine how the features of a dataset should split nodes to form the tree.**

5.Are unregularized decision-trees prone to overfitting? If yes, why?
**ANS :- Decision trees are prone to overfitting, especially when a tree is particularly deep. This is due to the amount of specificity we look at leading to smaller sample of events that meet the previous assumptions. This small sample could lead to unsound conclusions.**

6.What is an ensemble technique in machine learning?
**ANS :- Ensemble methods are techniques that aim at improving the accuracy of results in models by combining multiple models instead of using a single model. The combined models increase the accuracy of the results significantly. This has boosted the popularity of ensemble methods in machine learning.**

7.What is the difference between Bagging and Boosting techniques?
**ANS :- Bagging is a technique for reducing prediction variance by producing additional data for training from a dataset by combining repetitions with combinations to create multi-sets of the original data. Boosting is an iterative strategy for adjusting an observation's weight based on the previous classification.**

8.What is out-of-bag error in random forests?
**ANS :- The out-of-bag (OOB) error is the average error for each calculated using predictions from the trees that do not contain in their respective bootstrap sample. This allows the RandomForestClassifier to be fit and validated whilst being trained.**

9.What is K-fold cross-validation?
**ANS :- Cross-validation is a resampling procedure used to evaluate machine learning models on a limited data sample. The procedure has a single parameter called k that refers to the number of groups that a given data sample is to be split into. As such, the procedure is often called k-fold cross-validation.**

**ANS:- Cross-validation is a resampling procedure used to evaluate machine learning models on a limited data sample. The procedure has a single parameter called k that refers to the number of groups that a given data sample is to be split into. As such, the procedure is often called k-fold cross-validation.23-May-2018**

10.What is hyper parameter tuning in machine learning and why it is done?
**ANS :- In machine learning, hyperparameter optimization or tuning is the problem of choosing a set of optimal hyperparameters for a learning algorithm. A hyperparameter is a parameter whose value is used to control the learning process. By contrast, the values of other parameters (typically node weights) are learned.**

11.What issues can occur if we have a large learning rate in Gradient Descent?
**ANS :- In order for Gradient Descent to work, we must set the learning rate to an appropriate value. This parameter determines how fast or slow we will move towards the optimal weights. If the learning rate is very large we will skip the optimal solution.**

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

**ANS :- It can only be used to predict discrete functions. Hence, the dependent variable of Logistic Regression is bound to the discrete number set. It is very fast at classifying unknown records. Non-linear problems can't be solved with logistic regression because it has a linear decision surface.**

13.Differentiate between Adaboost and Gradient Boosting.
**ANS:-AdaBoost is the first designed boosting algorithm with a particular loss function. On the other hand, Gradient Boosting is a generic algorithm that assists in searching the approximate solutions to the additive modelling problem. This makes Gradient Boosting more flexible than AdaBoost.**

14.What is bias-variance trade off in machine learning?
**ANS:- In statistics and machine learning, the bias–variance tradeoff is the property of a model that the variance of the parameter estimated across samples can be reduced by increasing the bias in the estimated parameters.**

15.Give short description each of Linear, RBF, Polynomial kernels used in SVM.
**ANS:-**
**Linear kernal :- Linear Kernel is used when the data is Linearly separable, that is, it can be separated using a single Line. It is one of the most common kernels to be used. It is mostly used when there are a Large number of Features in a particular Data Set.**

**RBF Kernal :- In machine learning, the radial basis function kernel, or RBF kernel, is a popular kernel function used in various kernelized learning algorithms. In particular, it is commonly used in support vector machine classification**

**Polynomial Kernel :-**
**The polynomial kernel is a general representation of kernels with a degree of more than one. It's useful for image processing.**
**There are two types of this:**
**a.Homogenous Polynomial Kernel Function**
**K(xi,xj) = (xi.xj)d, where '.' is the dot product of both the numbers and d is the degree of the polynomial.**
**b.Inhomogeneous Polynomial Kernel Function**
**K(xi,xj) = (xi.xj + c)d where c is a constant.**