

# School of Computer Science Engineering and Technology

Course- BTech  
Course Code- 301  
Year- 2022  
Date- 21-02-2022

Type- Core  
Course Name-AIML  
Semester- Even  
Batch- 4<sup>th</sup> Sem (SPL)

## 6 - Lab Assignment No. 6.1

**Objective:** To Implement Random forest classifier

**Problem Statement:** Build a Decision Tree Classifier using Sklearn for predicting the number of shares in social networks (popularity).

**About Dataset:** This dataset summarizes a heterogeneous set of features about articles published by Mashable in a period of two years.

<b>Data Set Characteristics:</b>	Multivariate	<b>Number of Instances:</b>	39797	<b>Area:</b>	Business
<b>Attribute Characteristics:</b>	Integer, Real	<b>Number of Attributes:</b>	61	<b>Date Donated</b>	2015-05-31
<b>Associated Tasks:</b>	Classification, Regression	<b>Missing Values?</b>	N/A	<b>Number of Web Hits:</b>	355284

### Steps

1. **Dataset:** Download the dataset from the link  
<https://archive.ics.uci.edu/ml/datasets/Online+News+Popularity>
2. Remove unwanted features such as url, timedelta
3. Convert the categorical features into numerical using one hot encoding if any.
4. Perform range normalization on numerical features not in the range of 0 to 1.
5. Split the dataset into training and testing using 75-25 division
6. Build a Random Forest regression model using Sklearn with default parameters. Predict the target values in the testing set. Apply regression metrics and visualize the results as graphs.
7. Playing with Random Forest: Change the following parameters of the random forest and analyze their performance for training and testing using the evaluation measures.
  - a. n\_estimators
  - b. criterion{"mse", "mae"}
  - c. max\_depth
  - d. min\_samples\_split
  - e. bootstrap
  - f. n\_jobs
  - g. min\_samples\_leaf
  - h. max\_features
  - i. random\_state
  - j. max\_leaf\_nodes
8. Compare the performance of the Random Forest model with other regression models such as linear regression, polynomial regression, decision tree regression etc.

**Suggested Platform:** Python: Azure Notebook/Google Colab Notebook, packages such as numpy, nltk, regular expression package re.

# School of Computer Science Engineering and Technology

## **Additional** (Not a part of the evaluation)

\*\*\* You can try different strategies to see whether testing error comes down or not. Strategies can be different 1. removal of some features, 2. normalization methods, 3. Shuffling of training samples. Check the model error for the testing data for each setup.

\*\*\*Random Forest Classifier: Pick a Regression dataset of your choice and perform training testing similar to above. Play with model parameters and analyse the results using regression measures.