

## Assignment 2 - PySpark Data Preprocessing and Clustering

Submission Deadline: 20<sup>th</sup> April 2022

Submit the code and a Report in the form of a word document on google classroom.

**Question 1: Explore and preprocess the “Leaves” data given along with this assignment**

Load the given dataset into **Spark Data-Frames** and perform the following preprocessing on it.

- Handle missing values
- Identify if an attribute has outliers or noise
- Apply measures of the central tendency and dispersion to **analyze numeric attributes**. That is, compute the mean, median, mode, range, variance, correlation for the attribute. Don't just give values explain analyze them.
- Would you apply preprocessing techniques like discretization or normalization on any attribute? Explain your answer. If yes, then apply the technique and share the results.

**Question 2: CLUSTER THE DATA using the PySpark built-in K-means and bisecting K-means clustering algorithm (this is provided in the Spark ML Library).**

- Cluster the given dataSet using atmost **3-4** attributes to avoid curse of dimensionality. You can select the attribute based on the preprocessing.
- Cluster the dataset for different value of K

*Run your algorithm for various values of K and different values of convergence, show the results in your report.*

*Use measures such as SSE(the sum of square error), silhouette co-efficient, and NMI (normalized mutual Index) to analyze the clustering results.*

*See Scikit for more info on the above measures*

[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.normalized\\_mutual\\_info\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.normalized_mutual_info_score.html)