

# The University of Texas at Arlington

### MACHINE LEARNING (CSE 6363)

#### FINAL PROJECT

## Hierarchical Clustering for Seed Categorization

Name: Kuldip Rameshbhai Savaliya

**UTA ID: - 1001832000** 

#### **Project's Objective**

- Implement Hierarchical Clustering on the UCI seed dataset to divide it in groups and implement KNN to identify the species.
- Implement agglomerative clustering and KNN.

#### **Dataset link:**

https://archive.ics.uci.edu/ml/machine-learning-databases/00236/

#### **Column Descriptions:**

- 1.area A,
- 2. perimeter P,
- 3. compactness  $C = 4*pi*A/P^2$ ,
- 4. length of kernel,
- 5. width of kernel,
- 6. asymmetry coefficient
- 7. length of kernel groove.

#### Steps:

- Grouping the seed dataset into n clusters using hierarchical agglomerative clustering.
- Creating new features for the datapoints in the given dataset based on other clusters (using cluster's centroid) and added them into the existing features of the dataset. Generate new features based on distance.
- Similar procedure is followed for testing dataset.
- New Features are scaled using min max.
- Now, training dataset have extra feature values that is calculated based on similarity between data points.
- Implement KNN classifier on the modified dataset to identify the cluster ids.
- Performance of the model on unknown data is measured by the K-fold cross validation method. (Dataset is split into k subsets)
- Find the suitable number of clusters for different single linkages, complete linkage, and average linkage.
- Make Predictions on unlabeled data.

#### **Code Structure**

• Implement Hierarchical clustering (agglomerative clustering) model to divide the training dataset into n clusters.

```
 class AgglomerativeClustering:

         __init___ - Initialization
         fit_predict - fit data and predicting labels
         clustertolabels - for clusters to labels
         d_matrix - distance matrix
         update_d - update distance matrix according to linkage
         update_cluster - merge cluster
         distance - distance between two points
```

Implement KNN for classification of labels for unlabeled dataset.

```
def KNN_classification (sample, k, X, y):
 dist_cartesian – calculate cartesian distance
 lbl_classify – find label
 KNN classification - classification
```

 Implement k-fold cross validation method for estimating the performance of the model.

- Code Block for the value of clusters, neighbours, and linkage.
  - n\_clusters = [3,4,5,6,7]
    knn = [3,5,7,9,11]
    linkages = ['single', 'complete', 'average']
- Code block is used to make predictions on test data.
- Run KNN on given dataset (without adding new features) and find the accuracy.

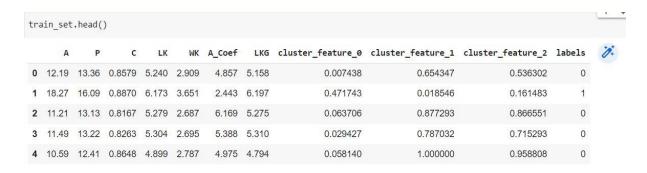
#### Results

- Experiments with different values of clusters, linkage, and neighbours
- n\_clusters = [3,4,5,6,7]
- knn = [3,5,7,9,11]
- linkages = ['single', 'complete', 'average']
  - result = k\_Fold (train\_set, 8, k)
- Accuracy of KNN is calculated based on mean of the accuracies of the folds found in list. (Steps repeated for different values of k)
- Average accuracy is calculate based on mean of the accuracies of different neighbours. (Steps repeated for different number of clusters)
- Best number of clusters find based on the highest accuracy for the linkage.
  - n\_clusters[np.argmax(Accuracy)]

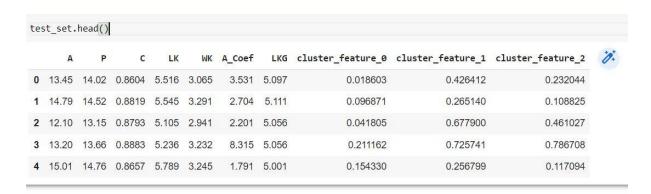
#### Dataset:

	Α	Р	C	LK	WK	A_Coef	LKG	target	%
0	15.26	14.84	0.8710	5.763	3.312	2.221	5.220	0	
1	14.88	14.57	0.8811	5.554	3.333	1.018	4.956	0	
2	14.29	14.09	0.9050	5.291	3.337	2.699	4.825	0	
3	13.84	13.94	0.8955	5.324	3.379	2.259	4.805	0	
4	16.14	14.99	0.9034	5.658	3.562	1.355	5.175	0	

#### Training data with features:



#### Testing data with features:



#### **Output:**

#### Scenario for single Linkage:

No of clusters:3

Accuracy: 97.4666666666668 %

No of clusters:4

No of clusters:5

Accuracy: 97.33333333333333 %

No of clusters:6

Accuracy: 94.80000000000001 %

No of clusters:7

Accuracy: 95.0666666666668 % Best scenario for no of clusters: 4

#### Scenario for complete Linkage:

No of clusters :3

Accuracy: 99.0666666666668 %

No of clusters :4 Accuracy : 97.2 % No of clusters :5

Accuracy: 90.93333333333334 %

No of clusters :6

Accuracy: 94.6666666666669 %

No of clusters:7

Accuracy: 91.73333333333335 % Best scenario for no of clusters: 3

#### Scenario for average Linkage:

No of clusters:3

No of clusters:4

Accuracy: 97.33333333333334 %

No of clusters:5

No of clusters :6

No of clusters :7

From above results, we can see that for single linkage best number of clusters 4 and for complete linkage best number of clusters is 3 and for average linkage best number of clusters is 4.

## Comparision between original dataset vs modified dataset:

K value	Modified dataset	Original dataset			
	(no. of cluster 3)				
3	0.9821	0.8666			
5	0.9761	0.8933			
7	0.9821	0.9066			
9	0.9761	0.9133			
11	0.9702	0.8933			

From above results, we can see that after adding new features in the dataset we get better accuracy compare to original dataset.