

PAM: introduction

The Partitioning Around Medoids (PAM) algorithm (also known as k-medoids) is a clustering technique aiming at minimizing the sum of dissimilarities between points in a dataset and the nearest medoid, where a medoid is the most centrally located point in a cluster. The main objective function is to minimize the total cost:

$$\sum_{i=1}^n \min_{j \in \{1, 2, \dots, k\}} d(x_i, m_j),$$

where $d(x_i, m_j)$ is the dissimilarity between point x_i and medoid m_j , and k is the number of clusters which has to be assessed beforehand. PAM is a better alternative than k-means when dealing with non-Euclidean distances or datasets with noise and outliers, as it is more robust to such irregularities.

'iris' dataset

iris: dataset of 150 observations x 5 variables.

Sepal.Length: continuous variable

Sepal.Width: continuous variable

Petal.Length: continuous variables

Petal.Width: continuous variable

Species.: discrete variables, setosa, virginica, versicolor

```
1 # load dataset
2 from sklearn.datasets import load_iris
3
4 # Load the iris dataset
5 iris = load_iris()
6 iris_df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
7 iris_df
8 sepal length (cm) sepal width (cm) petal length (cm) petal width (cm)
9 0 5.1 3.5 1.4 0.2
10 1 4.9 3.0 1.4 0.2
11 2 4.7 3.2 1.3 0.2
12 3 4.6 3.1 1.5 0.2
13 4 5.0 3.6 1.4 0.2
```

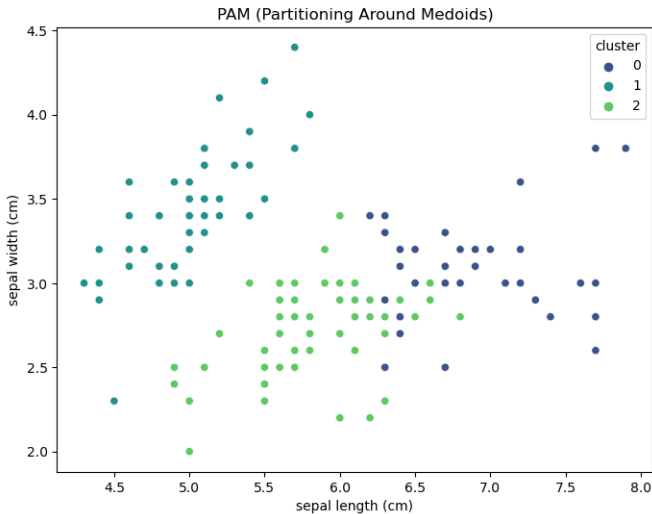
Summary

In order to perform a PAM, we call the function 'KMedoids()' from '*sklearn_extra.cluster*', containing three essential arguments: 'method', 'metric' and 'n_clusters'.

```
1 # Calculate pairwise distances
2 metric = 'manhattan'
3 dissimilarity_matrix = pairwise_distances(iris_df, metric='manhattan')
4
5 # Perform PAM clustering using KMedoids from sklearn_extra
6 pam = KMedoids(init='random', method='pam', metric='precomputed', n_clusters=3,
7               random_state=2024)
8 pam.fit(dissimilarity_matrix)
9
10 # Extract medoids and clustering results
11 medoids = pam.cluster_centers_
12 clustering_results = pam.labels_
13 clustering_results
```

Visualizing PAM clusters

On Iris dataset, using Manhattan distance metric



Main observations

- The PAM algorithm successfully grouped the Iris dataset into three clusters, with each cluster roughly corresponding to one of the true species.
- Some overlap is observed between the clusters, particularly between the clusters for versicolor and virginica, indicating difficulty in perfectly separating these two species.
- The comparison of true species labels with PAM clusters reveals that the algorithm captures the overall structure of the dataset but may misclassify some points, especially where species characteristics are similar.

References

<https://www.datanovia.com/en/lessons/k-medoids-in-r-algorithm-and-practical-examples/>

Python

<https://www.python.org/>