

Implementation of K-Means

Raúl Maximiliano Urrutia Hernández

November 7th, 2019

Algorithm

The algorithm that I used to make the implementation is (taken from: M. Zaki and W. Meira, *Data Mining and Analysis Fundamental Concepts and algorithms.*)

Algorithm 13.1: K-means Algorithm

K-MEANS (D, k, ϵ):

```
1  $t = 0$ 
2 Randomly initialize  $k$  centroids:  $\mu_1^t, \mu_2^t, \dots, \mu_k^t \in \mathbb{R}^d$ 
3 repeat
4    $t \leftarrow t + 1$ 
5   // Cluster Assignment Step
6   foreach  $x_j \in D$  do
7      $j^* \leftarrow \arg \min_i \{ \|x_j - \mu_i^t\|^2 \}$  // Assign  $x_j$  to closest centroid
8      $C_{j^*} \leftarrow C_{j^*} \cup \{x_j\}$ 
9   // Centroid Update Step
10  foreach  $i = 1$  to  $k$  do
11     $\mu_i^t \leftarrow \frac{1}{|C_i|} \sum_{x_j \in C_i} x_j$ 
12 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\| \leq \epsilon$ 
```

Implementation

First the algorithm randomly initializes the centers:

```
24         #Randomly initialize k centers
25         random_state = np.random.mtrand._rand
26         seeds = random_state.permutation(n_samples)[:self.k]
27         self.centers = X[seeds]
```

In the cluster assignment step, each point is assigned to the closer center:

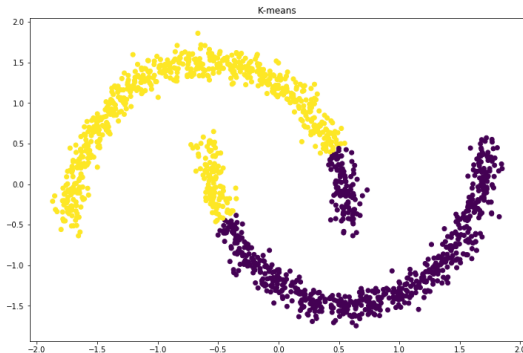
```
36         #Cluster assignment step
37         for entry in X :
38             distances = [np.linalg.norm(entry - center) ** 2
39                          for center in self.centers]
40             label = distances.index(min(distances))
41             groups[label].append(entry)
```

In the center update step, each center is updated to the average or mean of the cluster:

```
46         #Centers update step
47         for l in groups :
48             self.centers[l] = np.average(groups[l], axis = 0)
```

Testing the implementation

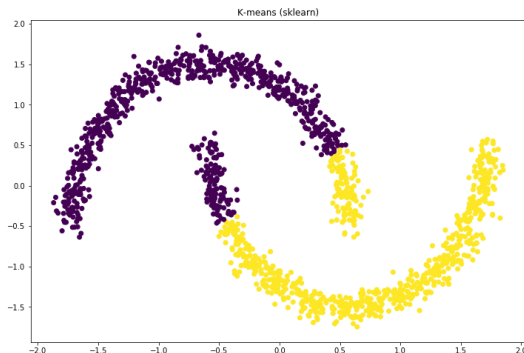
The scatter plot of the clustering of the make_moons dataset made by this implementation goes as follows



The adjusted Rand Index takes a value of 0.48407, and the adjusted Mutual information is 0.38492.

Testing the implementation

The scatter plot of the clustering made by sklearn kmeans

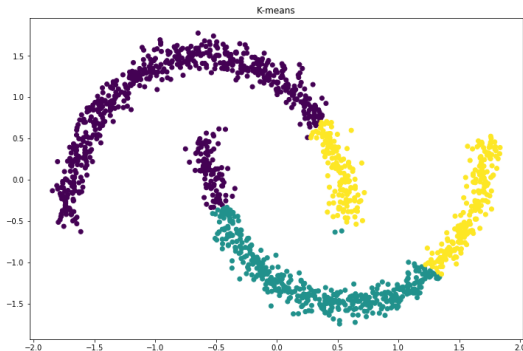


With adjusted Rand Index of 0.48407 and adjusted Mutual information of 0.38492.

As we can see from the plots and the metrics, both results (from the implementation and from the sklearn kmeans) are the same.

Testing the implementation

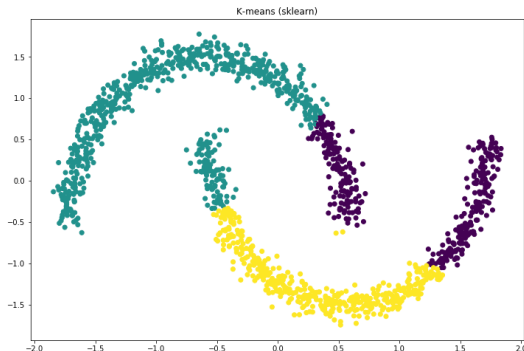
Lets make another test, now with 3 clusters. The scatter plot of the clustering made by the implementation is



This time, the adjusted Rand Index takes a value of 0.39336, and the adjusted Mutual information is 0.31165.

Testing the implementation

The scatter plot of the clustering made by sklearn kmeans is



With adjusted Rand Index of 0.38748 and adjusted Mutual information of 0.30811.

The metrics of both results are almost the same (slightly different due to the random nature of k-means).