

Homework 6

Software Quality

Group 06

Okikioluwa Ojo (100790236)

Date: Mar 5, 2024

GitHub Link: <https://github.com/okikio-school/Homework-6>

You will use the following open source java Machine Learning package to do some data clustering. <http://java-ml.sourceforge.net/>

You read the getting started section to set up this library. The assignment is for you to compare the performance of three clustering algorithms (including Kmeans and you choose any two other algorithms). Use the Iris flower dataset we discussed in class.

Note you will use the same method signature to do the clustering ... in spite of a different algorithm being employed? Can you tell why???

You need to explore how to evaluate each of the three algorithms. You need to read the documentation to explore how and what to evaluate? Did you use the same method signature to evaluate as well? Why?

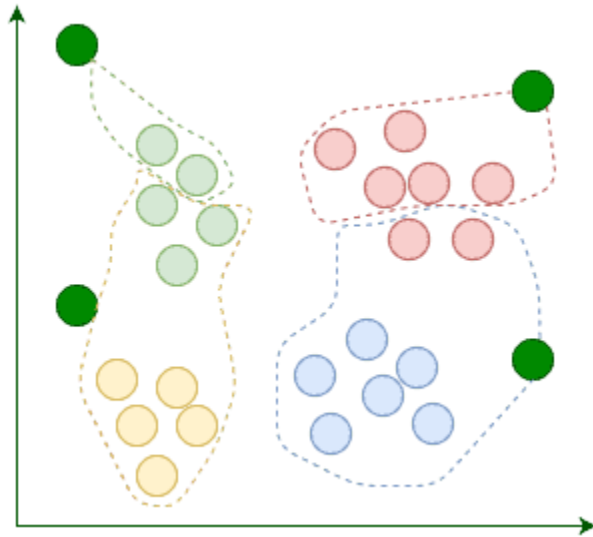
The clustering algorithm can be summarized as finding a pattern and reinforcing it into clusters, the basic shape of the algorithm is one that takes data and returns a modified set, where the data is placed into clusters.

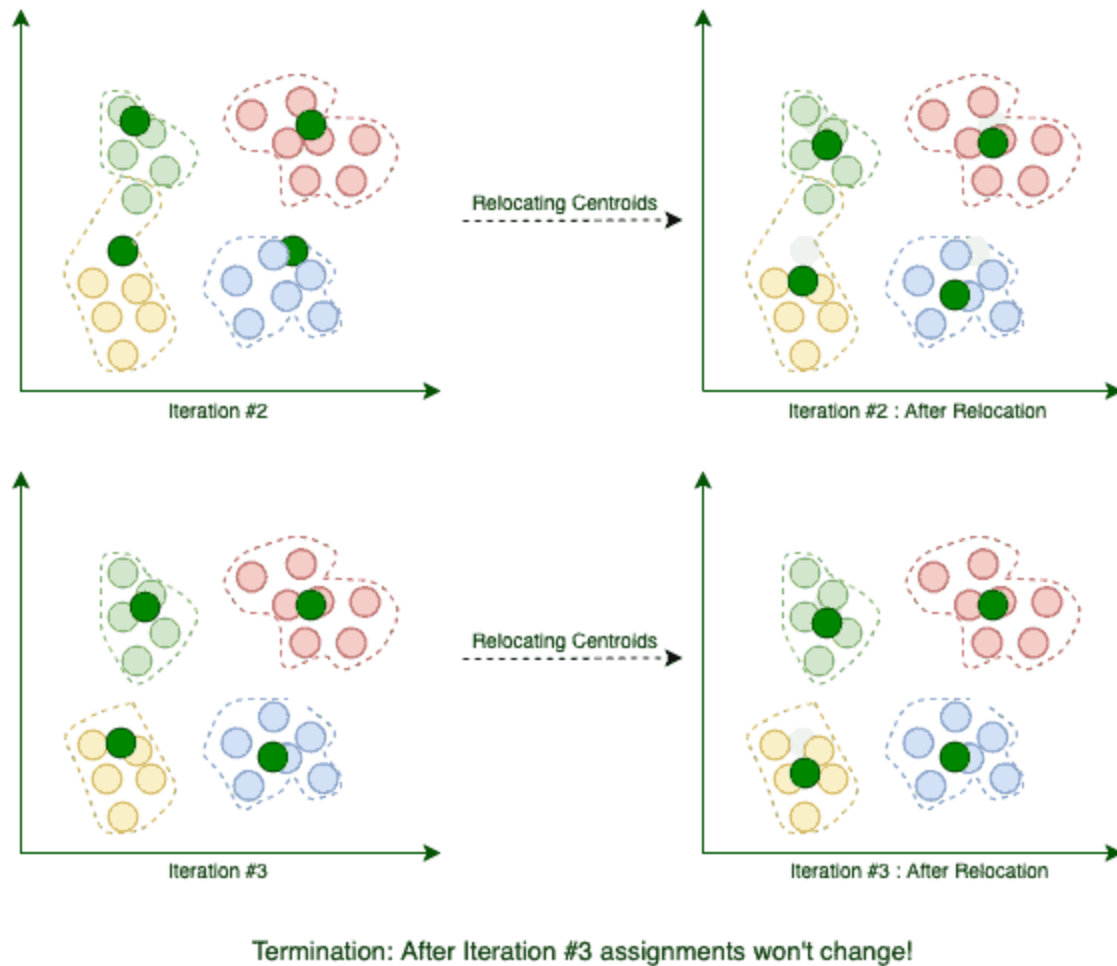
- Kmeans is a centroid-based clustering algorithm, it partitions datasets into a predefined set of clusters which is what we refer to when as `k` in KMeans (in the demo code I use 4 clusters, which is the default JavaML comes with). The KMeans algorithm iteratively assigns each data point (instance) to the nearest cluster center (this is the centroid part of the algorithm), it uses distance to identify the cluster center. The centroids are recalculated until they stabilize (e.g. the centroids no longer really change position so we classify the centroid as having reached stability).
 - KMeans aims to partition the data set into distinct non-overlapping clusters, where each instance belongs to the cluster with the nearest mean (nearest centroid), as KMeans minimizes the variance within a cluster. Basically we want more of the same type of data within the cluster.
 - KMeans starts by initializing K centroids randomly.
 - Each instance is assigned to the nearest centroid based on distance.
 - The centroids of the clusters are recalculated as the mean of all points assigned to the cluster.
 - This process repeats until the centroids no longer change significantly

To evaluate the performance of KMeans we can use metrics like Inertia, and the Silhouette Score. Inertia refers to the sum of squared distances to their closest cluster center. Lower inertia implies a better algorithm. For the Silhouette Score, it measures an instance to other instances in its own cluster. The silhouette

ranges from -1 to 1, a high value indicates that the object matches its own cluster fairly well, and implies that the other clusters are poor matches.

Wonderful resource <https://www.baeldung.com/java-k-means-clustering-algorithm>





- Farthest First, is a simpler implementation of KMeans, it selects centroids that are as far away from each other as possible, without iterating. What we mean by without iteration is that we do not repeatedly refine the choice of centroid through multiple passes of the data. Instead, we identify centroids based on the principle of maximizing the minimum distance between any set of points chosen to be centroids; this helps us spread out the centroids as much as possible within the dataset from the outset. By not using iteration we keep the algorithm simple, fast, and efficient.
 - The algorithm starts by selecting a random point as the first centroid.
 - For each subsequent centroid the algorithm selects the data point that is farthest away from any centroid that's already been chosen.
 - We continue this process until the desired number of centroids (k) have been met.
 - One key thing to note is that unlike KMeans, which aims to minimize the differences between instances in a cluster, FarthestFirst aims to have a

more diverse coverage for the dataset, all without needing to adjust the centroids based on cluster assignment.

Like KMeans, the Farthest First algorithm can be evaluated using inertia and silhouette score.

- Cobweb is an incremental clustering algorithm which uses a tree structure to dynamically organize data into hierarchical trees of clusters, generally meant for categorical data. When a new instance is introduced, the Cobweb algorithm either merges it into an existing cluster, creates a new cluster, splits an existing cluster or merges 2 clusters together. The decision is based on maximizing the probabilistic measure (category utility function).
 - Creates an empty tree.
 - For each data point (instance), the algorithm traverses the current tree to find the best cluster (each cluster counts as a node in the tree) that can accommodate the new instance, based on probabilistic measure (this is often called the category utility function) that evaluates how well the instance fits into a potential cluster.
 - Depending on the probability, Cobweb will either add the instance to a cluster node, create a new cluster if it doesn't fit into existing clusters, merge to clusters nodes into one if it improves the probabilistic measure, or split a cluster node into 2 separate clusters if it improves the probabilistic measure
 - This is then repeated for each instance, until the all the points have been defined under a cluster

The category utility function is a criterion for determining the quality of a clustering structure based on how well each of the instances fit into the category of the cluster node. The measure assesses both the probability of predicting attribute values within clusters (basically if the category is teens, the algorithm should be able to accurately predict that the age from 13 - 19, if it can't as in the category utility function is incorrect it takes one of the actions listed above ^). Cobweb, is very similar to KMeans in the results they're aiming for, but unlike KMeans which aims to minimize differences in a cluster (passive), Cobweb aims to maximize similarities in a cluster (active)

Evaluating Cobweb can be tricky, but one way of doing it is to measure the overall classification quality against known categories that exist in the data.

Can you print the output of each algorithm in a manner that shows the different clusters?

```
root@Okiki-PC → Homework 6 git:(main) X mvn exec:java

[INFO] Scanning for projects...
[INFO]
[INFO] -----< com.ontariotechu.sofe3980U:NextDate >-----
[INFO] Building AppTest 1.0.0
[INFO] from pom.xml
[INFO] -----[ jar ]-----
[INFO]
[INFO] --- exec:3.2.0:java (default-cli) @ NextDate ---
[INFO] Iris dataset loaded successfully - 150 instances

(KMeans) Cluster count: 4, Task duration: 6 milliseconds
Cluster Centers: [[{5.579999999999998, 2.6333333333333337, 3.986666666666667, 1.2333333333333334},null], {[5.005999999999999, 3.4180000000000006, 1.464, 0.24399999999999999];null}, {[7.086956521739132, 3.1260869565217386, 6.013043478260869, 2.143478260869565];null}, {[6.293617021276596, 2.8999999999999995, 4.951063829787233, 1.7297872340425529];null}]
Number of Points in Each Cluster: [30, 50, 23, 47]
Inertia: 57.47327326549492
Silhouette Score: 0.4951199955302146

(Farthest First) Cluster count: 4, Task duration: 0 milliseconds
Cluster Centers: [[{6.048148148148147, 2.8160493827160495, 4.6234567901234565, 1.567901234567901};null], {[4.818181818181818, 3.2212121212121207, 1.436363636363636, 0.2272727272727273];null}, {[7.173684210526318, 3.110526315789474, 6.110526315789474, 2.136842105263158];null}, {[5.370588235294117, 3.8, 1.5176470588235293, 0.2764705882352942];null}]
Number of Points in Each Cluster: [81, 33, 19, 17]
Inertia: 87.90682427995128
Silhouette Score: 0.38375508476125586

(Cobweb) Cluster count: 11, Task duration: 19 milliseconds

[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 0.988 s
[INFO] Finished at: 2024-03-11T05:36:26-04:00
[INFO] -----
root@Okiki-PC → Homework 6 git:(main) X
```