

Assignment 4: K-Means Clustering

Assignment Overview

In this assignment you will implement the k-means clustering algorithm. You will use it to cluster a data set. For this assignment, we will use a data set from the UC Irvine Machine Learning Repository at:

<https://archive.ics.uci.edu/ml/index.html>.

Write your own code!

For this assignment to be an effective learning experience, you must write your own code! **Do no share code with other students in the class!!**

Here's why:

- The most obvious reason is that it will be a huge temptation to cheat: if you include code written by anyone else in your solution to the assignment, you will be cheating. As mentioned in the syllabus, this is a very serious offense, and may lead to you failing the class.
- However, even if you do not directly include any code you look at in your solution, it surely will influence your coding. Put another way, it will short-circuit the process of you figuring out how to solve the problem, and will thus decrease how much you learn.

So, just don't look on the web for any code relevant to this problem. Don't do it.

Format of data file

The data file that you are clustering is a database related to car evaluation. A complete description can be found here:

<https://archive.ics.uci.edu/ml/machine-learning-databases/car/car.names>

The input data is provided named as 'input_car'.

Each line of the file looks something like this: vhigh,vhigh,4,2,big,med,unacc

It consists of six attributes values and a text label.

The six attributes correspond to:

1. buying: buying price
2. maint: price of the maintenance
3. doors: number of doors of the car
4. persons: capacity in terms of persons to carry
5. lug_boot: the size of luggage
6. safety: estimated safety of the car

The last column is the class name, there are four classes:

- unacc: unaccurate
- acc: accurate
- good: good
- vgood: very good

Submission Details

What you will turn in: `firstname_lastname_clustering.py`

Program Description

Your program should do the following:

1. Read the data from the file. Use only the **floating point values for the clustering**. Don't discard the class information. While you can't use it for clustering, you will need it later for assigning names to the clusters and for checking the accuracy of the clusters.
2. **Apply the k-means algorithm to find clusters.** (There are 3 natural clusters in the case of the iris data.) (See below for more information on k-means.) **Use Euclidean distance** as your distance measure.
3. **Assign each final cluster a name** by choosing the most frequently occurring class label of the examples in the cluster. If the frequencies of two labels are the same, choose one randomly.
4. **Find the number of data points that were put in clusters in which they didn't belong** (based on having a different class label than the cluster name).

k-means algorithm:

Given k initial points that will act as the centroids of the clusters for the first iteration, you will run the standard k-means clustering algorithm that we discussed in class.

- For each point, place it in the cluster whose current centroid it is nearest
- After all points are assigned, update the locations of centroids of the k clusters
- Repeat for the specified number of iterations.

Output of your program

The program will produce output of the form: (If your output file is not exactly the same as the output sample, 20% points will be subtracted.)

Cluster <clustername1>:

(List of points in that cluster, one per line)

Cluster <clustername2>:

(List of points in that cluster, one per line)

Cluster <clustername3>:

(List of points in that cluster, one per line)

Number of points assigned to wrong cluster:

(number of points)

Running your code

```
python lastname_firstname_clustering.py dataFileName initialPoints k iter
```

where:

dataFileName is a string indicating the name of the data file to be clustered

initialPoints is a string indicating the name of a file that contains a list of data points that are to be used as the starting centroids for each cluster

k is an integer representing the number of clusters (three in the case of the iris data set)
iter is the number of iterations for the k-means clustering to run

The name of your python file is lastname_firstname_clustering.py, all in lowercase. Please follow the naming convention. Or you will lose 20% points.

You are required to set four parameters for the command line. Please don't set extra parameters and keep the order of the parameters exactly same as the sample command line. Or you will lose 20% points

Testing your code

The sample command to execute is :

```
python clustering.py input_car initialPoints 4 10
```

Sample output file can be found in the homework holder. Please follow the format of the sample output file, or you will lose 20% points, there is no sequence for the lines in each cluster.