

HOMEWORK #2

Software Project (0368-2161)

Due Date: 15/06/2021

1 Introduction

The K-means++ algorithm is used to choose initial centroids for the K-means algorithm. In this assignment you will implement this algorithm in Numpy and integrate it with the K-means algorithm of HW1 that will be ported to a C extension using the C API.

The goals of the assignment are:

- Port your existing C code into a C extension using the C API.
- Experience the Numpy, Pandas and other external packages.

This assignment relies heavily on the previous one and should be read with that in mind. As in HW1, this assignment will be tested and graded on Nova.

1.1 K-means++

Given a set of N datapoints $x_1, x_2, \dots, x_N \in R^d$, the goal is to find K initial centroids $\mu_1, \mu_2, \dots, \mu_K \in R^d$ where $K < N$.

This algorithm will replace step 1 of the K-means algorithm presented in HW1.

Algorithm 1 k-means++

Input: N datapoints $x_1, x_2, \dots, x_N \in R^d$, $K < N$

Output: K initial centroids, $\mu_1, \mu_2, \dots, \mu_K \in R^d$

1: Select μ_1 randomly from x_1, x_2, \dots, x_N

2: $Z = 1$

3: **repeat**

4: For each x_i

$$D_i = \min (x_i - \mu_j)^2, \forall j \ 1 \leq j \leq Z$$

5: $Z++$

6: Let μ_Z be randomly selected from the points x_1, x_2, \dots, x_N , with the probability of selecting x_i :

$$P(x_i) = \frac{D_i}{\sum_{i=1}^N D_i}$$

7: **until** $Z = K$

The algorithm starts by choosing a random observation as the first centroid. If $K=1$ the algorithm terminates. In step 3 we keep growing our centroids by randomly choosing observations with a weighted probability, favoring a distant observation from the centroids we already chose. **The way we calculate the probabilities ensures that observations we previously chose as centroids will have the probability 0 (zero) of getting reselected.** Note that all the centroids we chose are “as-is” observations. We have implemented a sampling method. Do not confuse these initial centroids with the calculated ones that represent the means of the clusters in the K-means algorithm.

2 Assignment Description

Implement the following files:

1. `kmeans_pp.py`: The “glue” and the main entry point of your assignment. It will contain all of the command-line argument interface, reading the data, the Numpy K-means++ implementation, the interface with your C extension and outputting the results.
2. `kmeans.c`: A C extension containing your K-means implementation from HW1 with step 1 of the algorithm (finding the initial centroids) replaced by values you’ll pass from the K-means++ algorithm implemented in `kmeans_pp.py`.
3. `setup.py`: The setup file.

2.1 Main Program (`kmeans_pp.py`)

This is the interface of the program, with the following requirements:

1. Reading user CMD arguments:
 - `k`: Number of required clusters.
 - `max_iter`: (Optional) argument determines the number of K-means iterations, if not provided the default value is 300.
 - `file_name_1`: The path to the file 1 that will contain N observations, the file extension is `.txt` or `.csv`.
 - `file_name_2`: The path to the file 2 that will contain N observations, the file extension is `.txt` or `.csv`.
2. Combine both input files by **inner join** using the first column in each file as a key.
3. Implementation of the k-means++ algorithm as detailed in [1.1](#).
4. Interfacing with your C extension:
 - (a) Import C module `mykmeanssp`
 - (b) Call the `fit()` method with passing the initial centroids, the datapoints and other arguments if needed.
 - (c) Get the final centroids that returned by `fit()`.
5. Outputting the following:

- The first line will be the indices of the observations chosen by the K-means++ algorithm as the initial centroids. We refer to the first observation index as 0, the second as 1 and so on, up until N - 1. Observation's index is given by the first column in each input file.
- The second line onwards will be the calculated final centroids from the K-means algorithm, separated by a comma, such that each centroid is in a line of its own as in HW1, apply `np.round()` with 4 decimal places on your output.

An example output (assuming $K = 3$, `max_iter = 100` and that the initial centroids resulting from the K-means++ algorithm are the 1st, 5th and 22nd observations):

```
$python3 kmeans_pp.py 3 100 input_1.txt input_2.txt
0,4,22
-4.2435,9.1568,5.4105,9.6870,-5.7564,-7.2314
3.3226,-1.3896,-9.1927,-6.0907,-0.9954,-8.7412
8.2239,-8.5714,-8.4985,0.8969,-8.2158,-2.3753
```

2.2 K-means Algorithm (`kmeans.c`)

In this file you will define your C extension which will be, mainly, your implementation of the K-means algorithm from HW1, excluding step 1 which will be replaced by the K-means++ algorithm result. Requirements of this extension:

1. The C extension module should be named `mykmeanssp`.
2. The module API provides a function called `fit()`:
 - (a) The Function receives the initial centroids, datapoints and other necessary arguments.
 - (b) Skip step 1 from HW1 and run the algorithm from HW1.
 - (c) Return the centroids.

You can use or implement any auxiliary functions within the module, but only expose to the API the one described above. It's up to you to decide which and what type of arguments you want to pass when calling the API, same thing about the return.

2.3 Setup (`setup.py`)

This is the build used to create the `*.so` file that will allow `kmeans_pp.py` to import `mykmeanssp`.

2.4 Build and Running

1. The extension must built cleanly (no errors, no warnings) when running the following command:

```
$python setup.py build_ext --inplace
```

2. After succcessfull build, the program must run as detailed in example [2.1](#).

2.5 Assumptions and requirements:

Note that the following list applies to all of the files in this assignment:

1. You may assume that the input file is in the correct format and that it is supplied.
2. Validate that the command line arguments are in correct format.
3. In the C implementation, you should log and terminate the program on any error in allocating memory.
4. You may import external includes (in C) or modules (in Python) that are not mentioned in this document.
5. Three input files and their corresponding output files examples will be provided within the assignment in Moodle.
6. Comment your code. If you use code from the internet, please indicate the source.
7. Make your code (C and Python) as efficient as possible, you'll notice that the Python's runtime becomes increasingly slow as the number of observations, clusters and dimensions grow, to measure the runtime, use the "time" command in Linux. **There is no requirement to output your runtimes.**

2.6 Submission

1. Please submit a file named `id1_id2_assignment2.zip` via Moodle, where `id1` and `id2` are the ids of the partners.
 - (a) In case of individual submission, `id2` must be 111111111
 - (b) Create the zip file using **only** linux cmd:

```
$zip -r id1_id2_assignment2.zip your_directory_name
```
2. The zipped file must contain **only** the following files (at the root, no other folders):
 - (a) `kmeans_pp.py`
 - (b) `kmeans.c`
 - (c) `setup.py`
 - (d) `bonus.py` (optional)

2.7 Remarks

For any question regarding the assignment, please post at the HW_2 discussion forum.

3 Optional Bonus (5 points)

This section covers an optional bonus. The goal of this bonus is to demonstrate the use of the elbow method to determine the optimal number of clusters for the k-means clustering.

3.0.1 Elbow Method

We will define the quality of the k-means clustering algorithm as inertia and define it to be the sum of squared distances of samples to their closest cluster center:

$$inertia = \sum_{i=1}^N (x_i - \mu_{x_i})^2$$

where μ_{x_i} is the centroid of the cluster x_i is in.

The idea of the elbow method is to run k-means clustering with a range of values of k and for each value of k calculate the inertia. Then, plot a line chart of the inertia for each corresponding value of k. If the line chart looks like an arm, then the "elbow" on the arm is the value of k that is the best. That elbow could be defined as a range of k's if it is unclear exactly where the elbow is at. We would like you to plot that line chart on the iris dataset, using the module sklearn . The iris dataset contains a 4-dimensional 150 observations. Use the load_iris() API to access the iris dataset. Run the sklearn k-means algorithm (using the k-means++ initialization and with a random_state=0) for the values of k ranging from k=1 till k=10 and plot the inertia for each value of k using the matplotlib module.

It will be submitted with the following requirements:

1. As a Python file named `bonus.py`
2. This program will only produce an output `elbow.png` in the program folder.
3. Doesn't require any input.
4. Has no tester file to check against.
5. Annotate the chosen k on the plot, where the 'elbow' is. See example below:

