

Name: Tan Ngai Chuan Alvin

AML Practicum 1: K-means Algorithm

Abstract

To implement k-means algorithm to learn the clustering partitioning method from the training data. To apply k-means group dots representing lightning strokes. To familiarize with determining the parameter for k-means.

Using $M = [0,0]$, we computed the error as 86.62495736289799.
Using first 5 XY-coords of X as centroids (K=5), the error is reduced to 9.73227985727934
After running calcMeans() and Group() (K=5), the error is reduced to 4.006100077883563
After running K=50, the error is reduced to 0.9216408380317802
After running K=100, the error is reduced to 0.5725156609427651

Given a choice of either K=50 and K=100, we would prefer the K=50 over K=100 as the difference between the 0.92 error of K=50 and 0.57 error of K=100 is less than 0.50. In fact, both error yield an error of less than 1.0. Coupled with the visual challenge of 'over-clustering' on a plot, we would prefer the K=50 over K=100.

Introduction

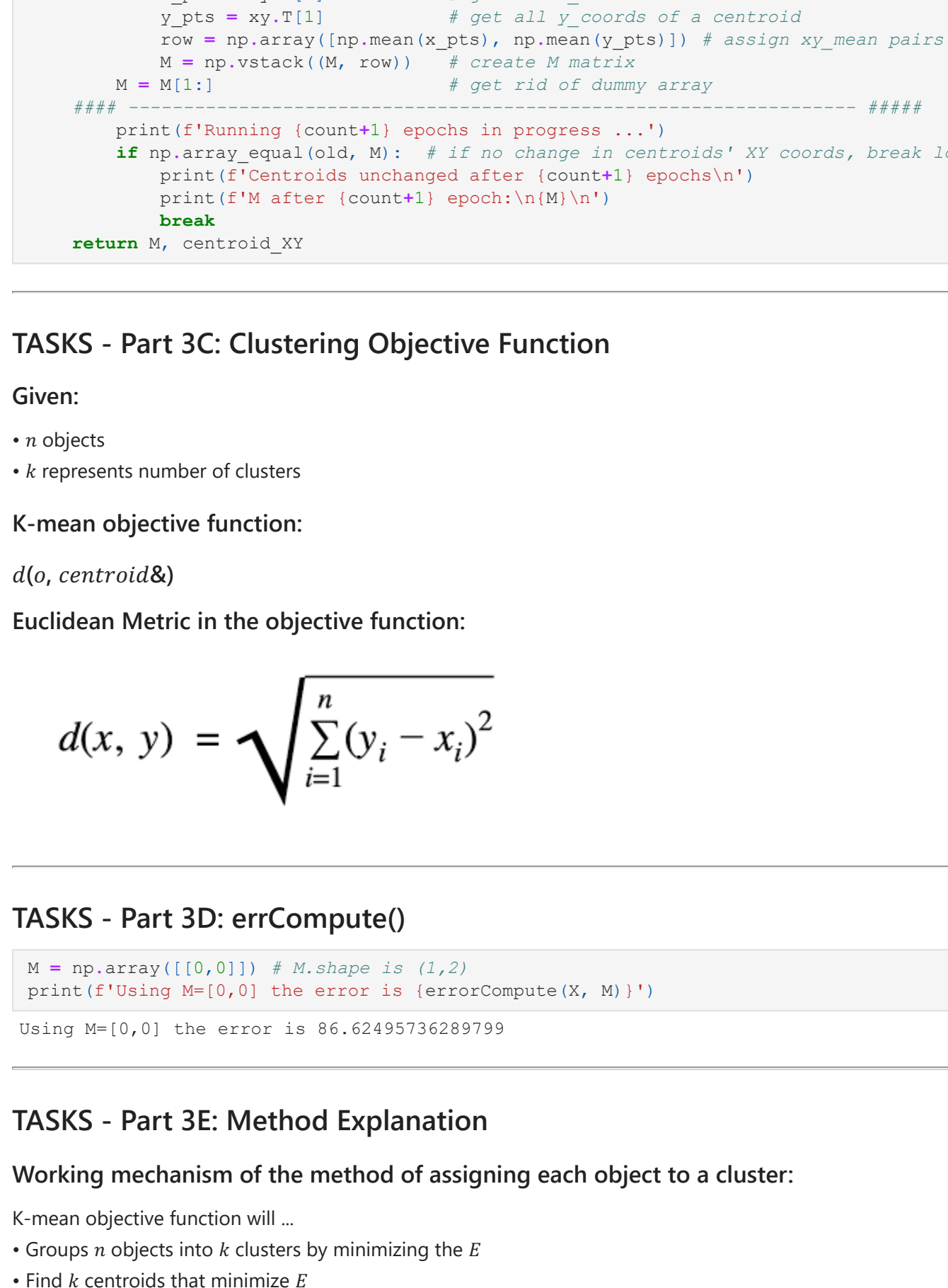
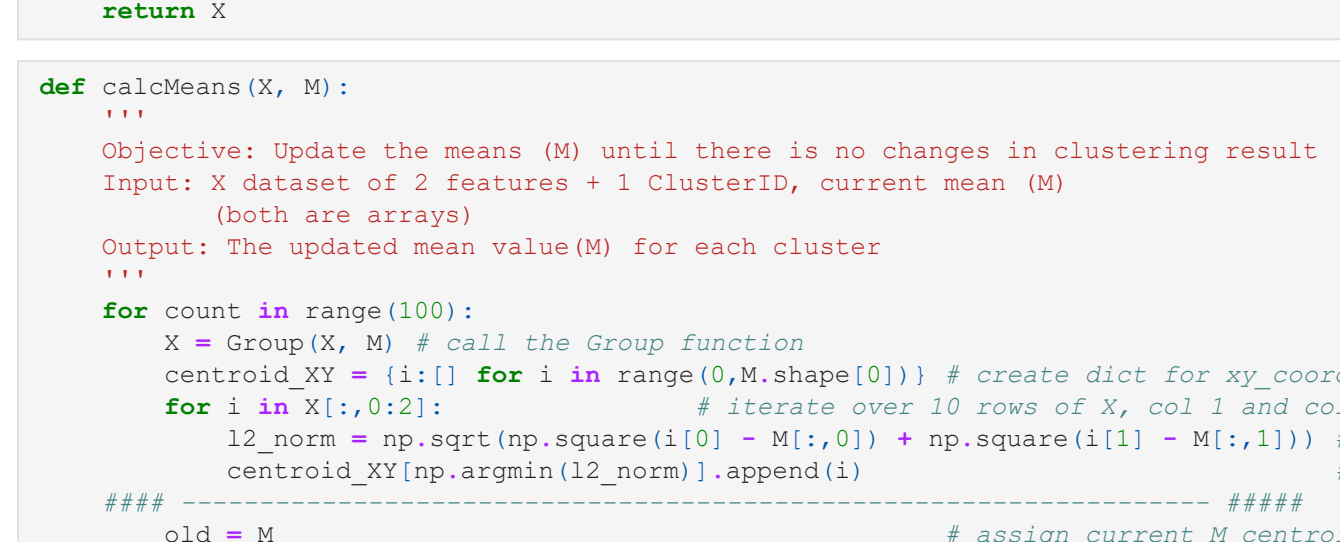
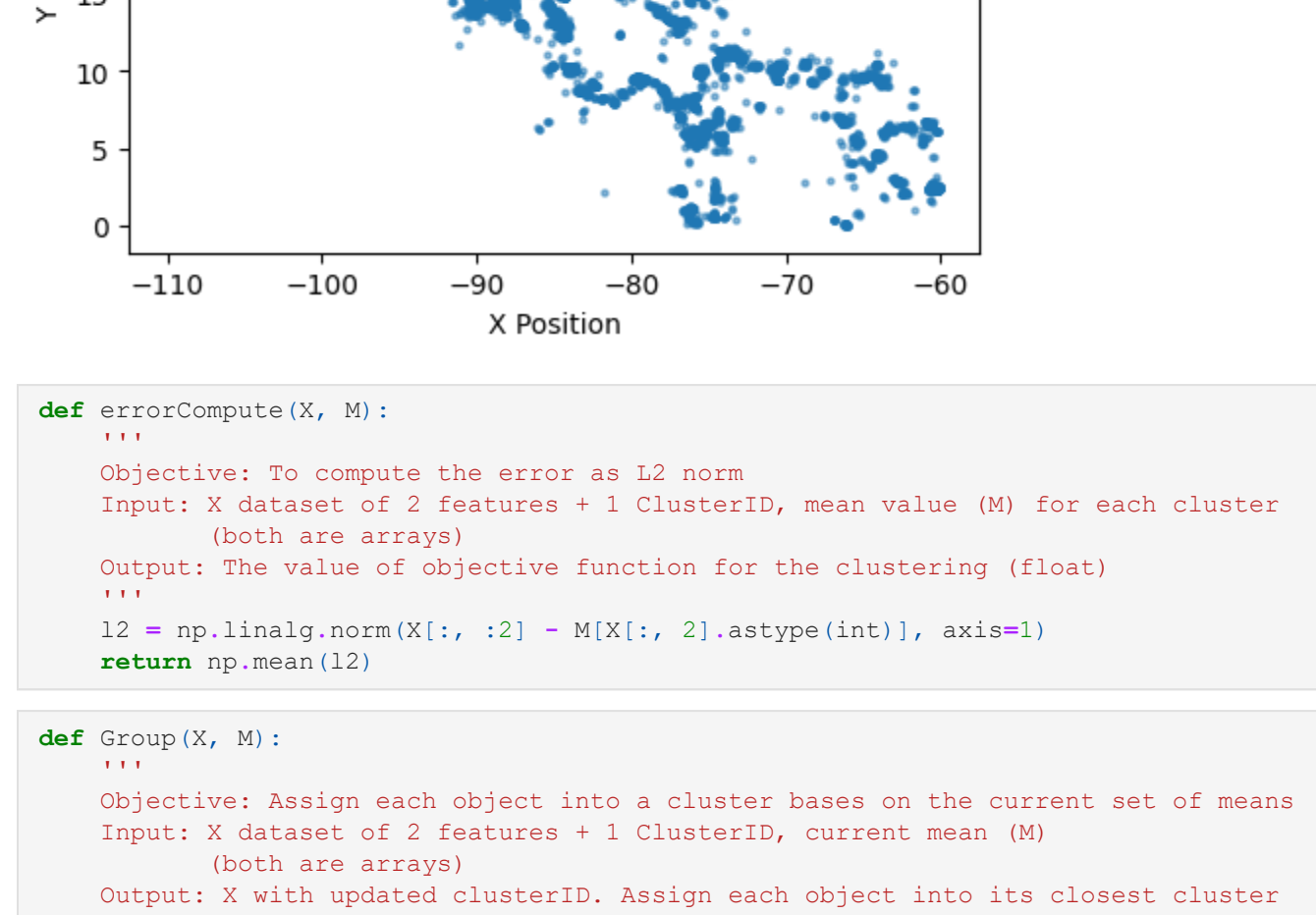
Lightning activities detect individual lightning discharge events all over the world. Understanding the lightning networks over every unique geography location over space and time will help us to better manage any economic damages and human safety that could arise from lightning activities.

For each lightning event, the data is stored as (time, amount of discharge, latitude, and longitude). It shows that when and where, there is a lightning stroke with certain amount of discharge.

We will implement a k-means clustering algorithm to help the meteorologists to group the lightning strokes basing on the position.

K-means clustering is one of the simplest and popular unsupervised machine learning algorithms. Typically, unsupervised algorithms make inferences from datasets using only input vectors without referring to known, or labelled, outcomes.

A cluster refers to a collection of data points aggregated together because of certain similarities. We define a target number k, which refers to the number of centroids needed in the dataset. A centroid is the imaginary or real location representing the center of the cluster. Every data point is allocated to each of the clusters through reducing the in-cluster euclidean distance. In other words, the K-means algorithm identifies k number of centroids, and then allocates every data point to the nearest cluster, while keeping the centroids as small as possible. The 'means' in the K-means refers to averaging of the data; that is, finding the centroid.



TASKS - Part 3C: Clustering Objective Function

Given:

- n objects
- k represents number of clusters

K-mean objective function:

$d(o, centroid_k)$

Euclidean Metric in the objective function:

$$d(x, y) = \sqrt{\sum_{i=1}^n (y_i - x_i)^2}$$

TASKS - Part 3D: errCompute()



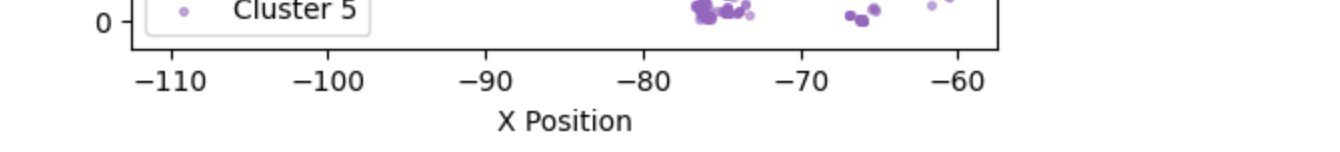
TASKS - Part 3E: Method Explanation

Working mechanism of the method of assigning each object to a cluster:

K-mean objective function will ...

- Groups n objects into k clusters by minimizing the E
- Find k centroids that minimize E
- Centroid is an actual object centrally located in a cluster

TASKS - Part 3F: Group()

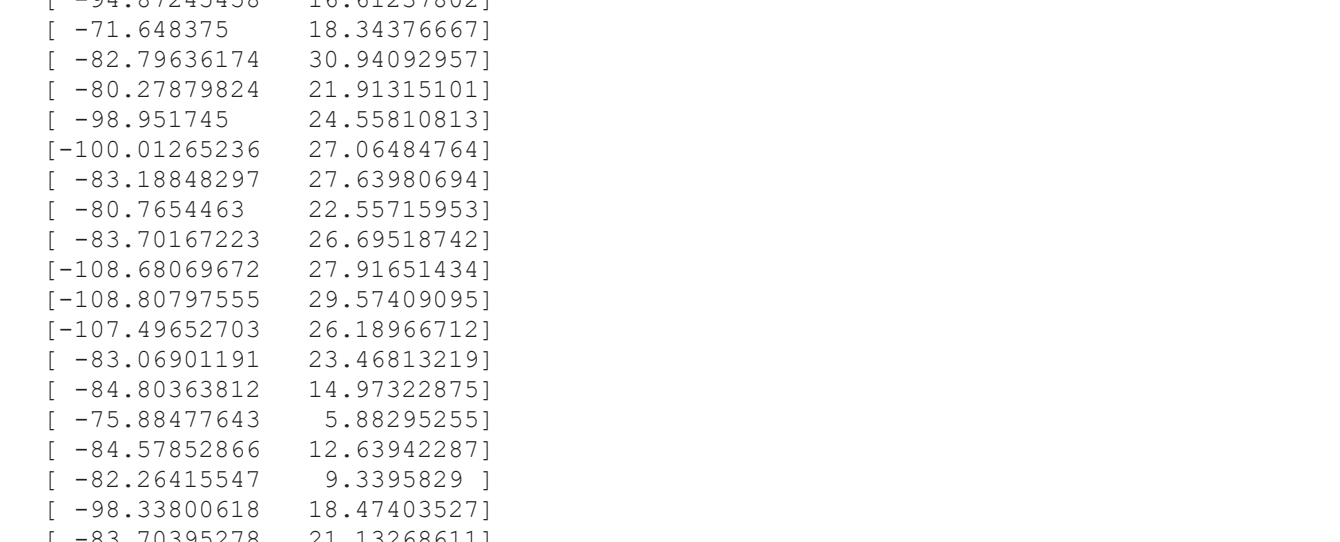
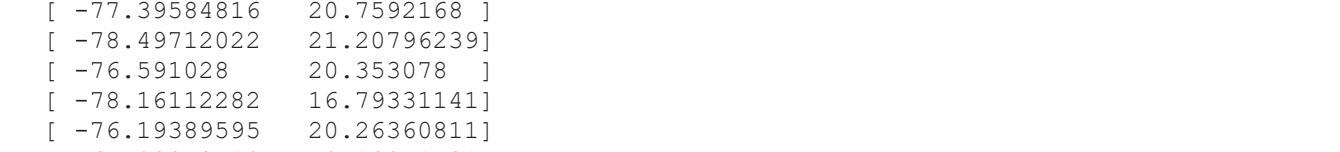
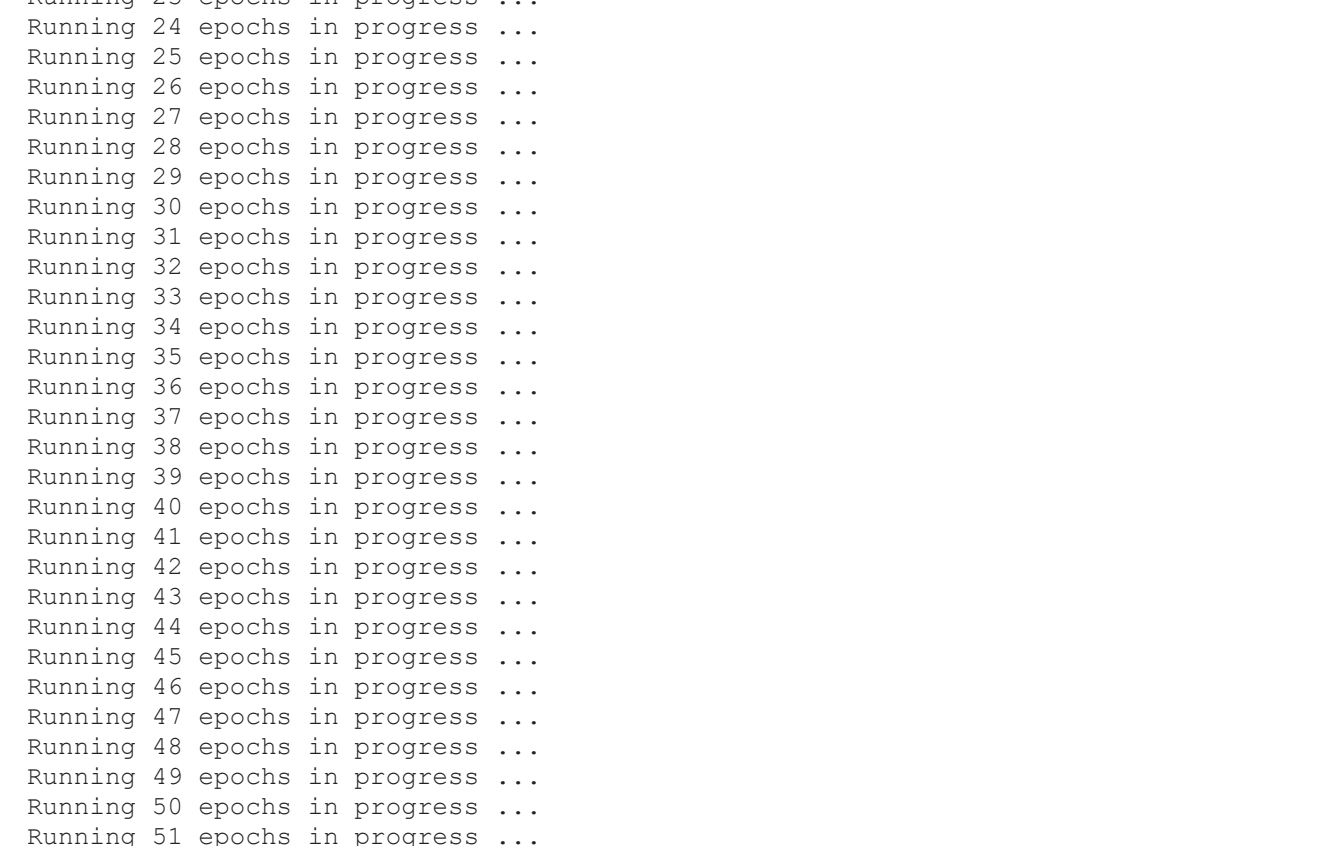


TASKS - Part 3G: Describe the method to compute new means

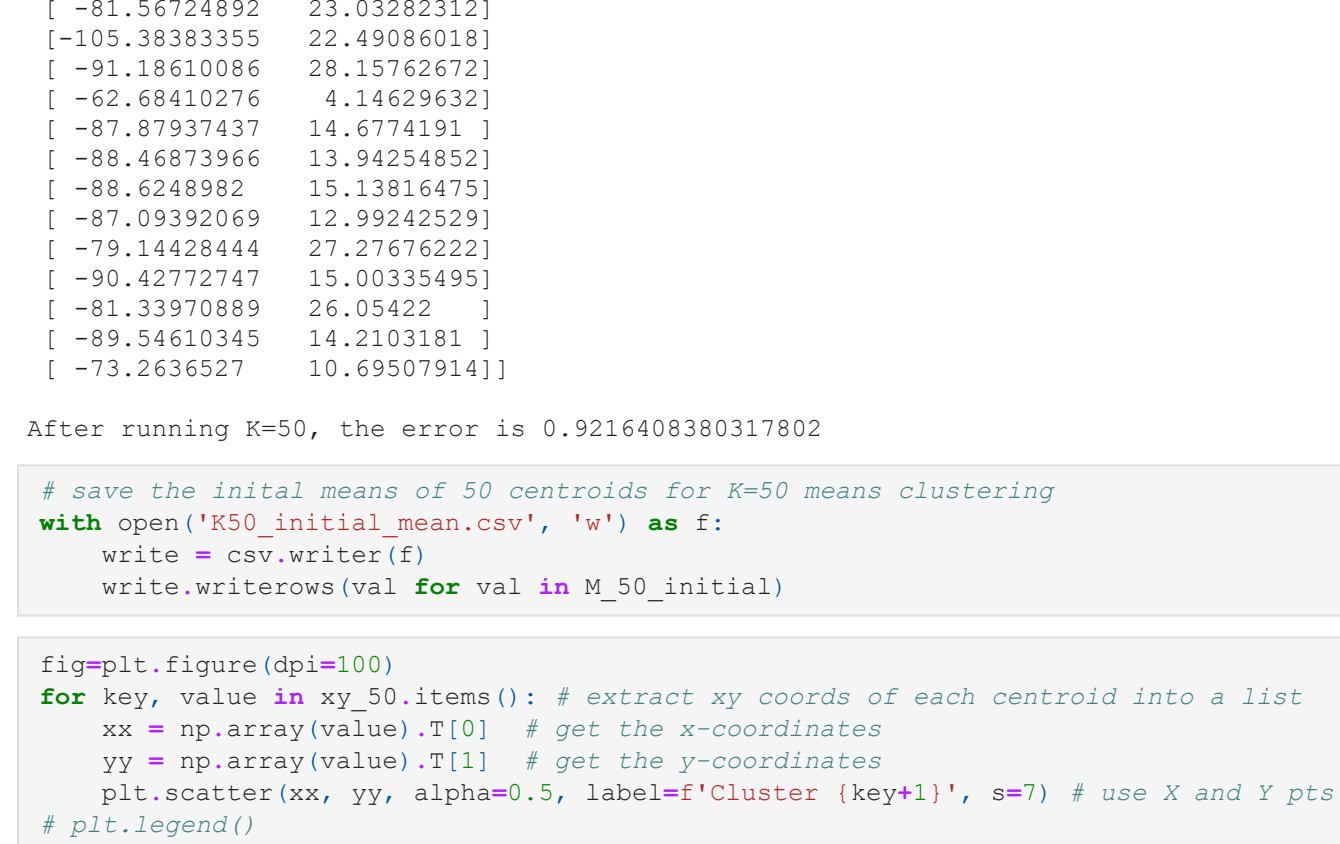
Steps to calculate the new means based on current clustering:

- Step1: Randomly select k number of centroids
- Step2: Assign each sample into the class represented by the closest centroids
- Step3: Update centroids as mean of the cluster based on the latest set of assigned data points
- Step4: Repeat step2 and step3 until convergence

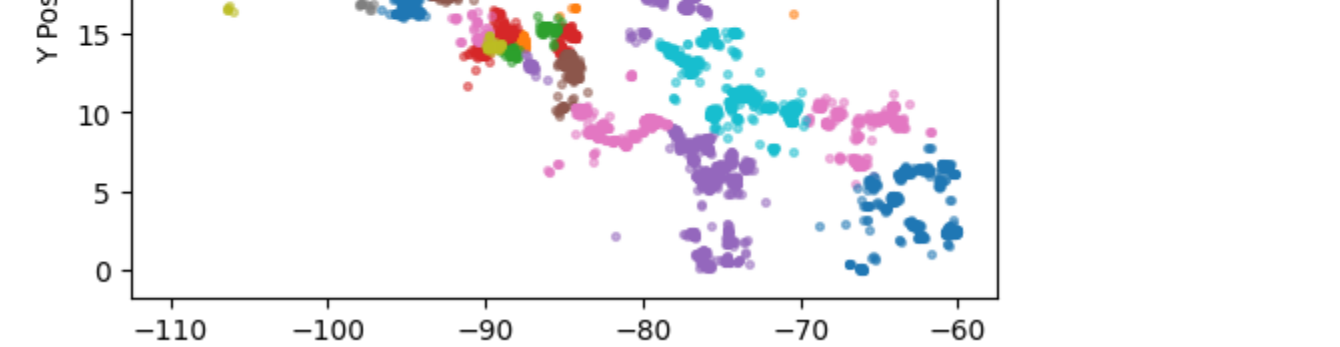
TASKS - Part 3H: Run k-means with K=5



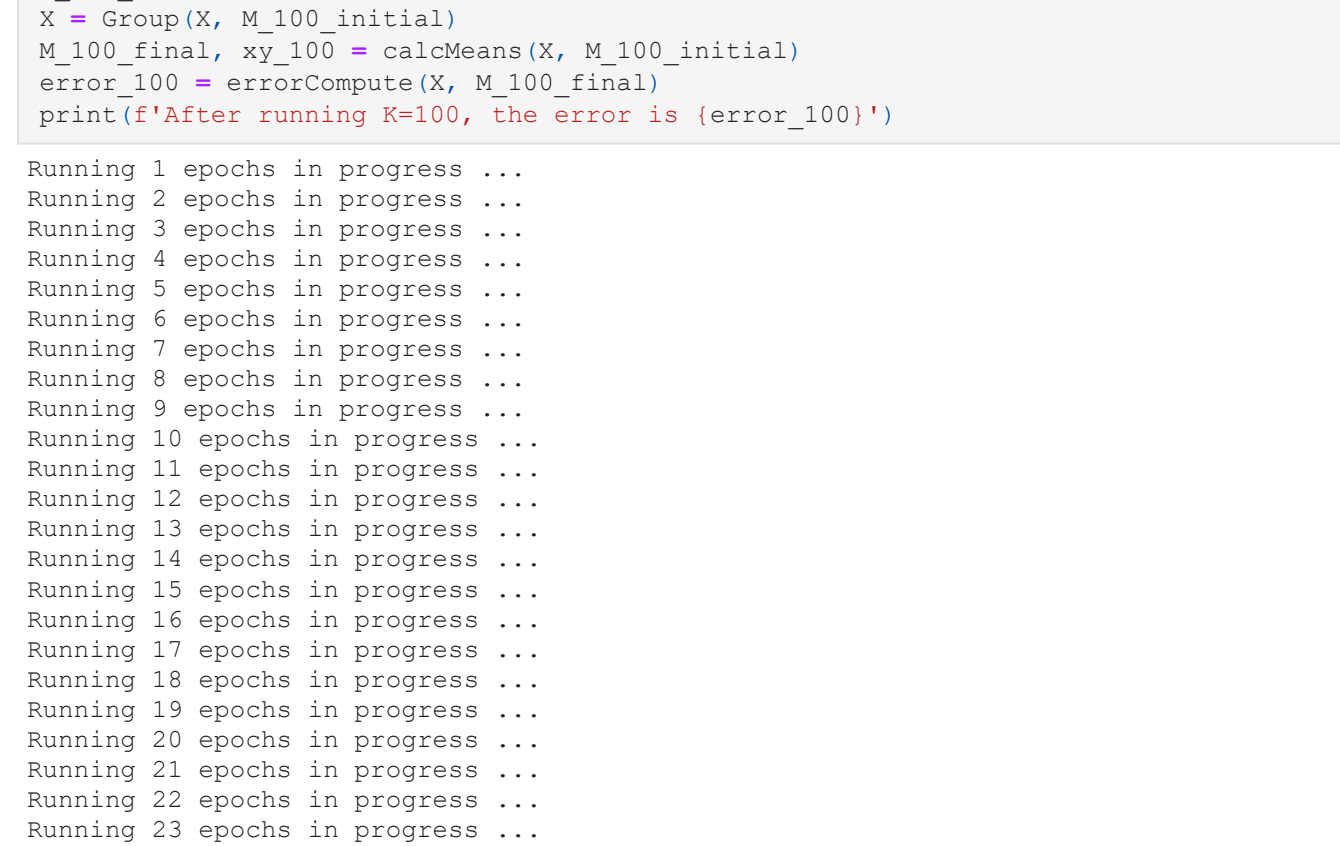
TASKS - Part 3I: Run k-means with K=50



After running K=50, the error is 0.9216408380317802



TASKS - Part 3J: Run k-means with K=100



After running K=100, the error is 0.5725156609427651



Discussion on which one (K=50 or K=100) is better

Based on the clustering plots of K=50 and K=100, we can see that both have at least 10 and 20 times more clusters than when K=5 respectively.

However, given a choice of either K=50 and K=100, we would prefer the K=50 over K=100 as the difference between the 0.92 error of K=50 and 0.57 error of K=100 is less than 0.50. In fact, both error yield an error of less than 1.0. Coupled with the visual challenge of 'over-clustering' on a plot, we would prefer the K=50 over K=100.

Clustering is useful and effective if we can find a right balance of error and visual-challenge of 'over-clustering', and given a choice, we would prefer K=5 for simple and easy labelling of lightning clusters.