

Exercise 7

A simple example of Spark Machine Learning

Prior Knowledge

Previous exercises on Spark

Learning Objectives

Understanding how to use Spark MLlib k-means clustering
Simple machine learning algorithms

Software Requirements

(see separate document for installation of these)

- Apache Spark 2.3.0
- Python 3.7
- Nano text editor or other text editor

Overall plan – finding geographic clusters of SFPD incidents

In a previous exercise, we managed to get all the data together for incidents in SF for 2014. This data included the [lat,long] for each incident.

In this exercise, we will simply provide that data to the K-means clustering algorithm that is built into Apache Spark MLlib.

In order to call the kmeans learning phase, we need to create an RDD where each element is a **numpy.array**.

This array will contain the [Y,X] data from the SFPD CSV file.

1. Load in the sfpd.csv file into a spark rdd. Map the rdd so that you just have the longitudes and latitudes of the incidents. Filter out any strange looking latitudes/longitudes (obviously not in San Francisco).
2. If you have an RDD with the Y, X values in a tuple like this (y,x), you can create the array using the following lambda:

```
from numpy import array
geoarray = rdd.map(lambda r: array([r[1],r[0]]))
```

3. Once you have an RDD of this form you can simply train the KMeans model like this:

```
from pyspark.mllib.clustering import KMeans, KMeansModel

numclusters = 5

clusters = KMeans.train(arrays, numclusters, maxIterations=10,
initializationMode="random")
```

4. This code will print out the cluster center points in lat,long format:

```
for arr in clusters.centers:
    clusterlist = arr.tolist()
    print("{}{}".format(clusterlist[0],clusterlist[1]))
```

5. Put all this together into a Spark program and run it. The output should be a set of data lines like:
37.7855185124,-122.408836684

You can cut and paste this data into the following webpage which will plot those points on a map:

<http://www.hamstermap.com/quickmap.php>

You can compare the results to the heat map of crime produced by Trulia here:

<http://www.trulia.com/blog/trends/trulia-local/>

Note that the clusters represent all incidents not just crime, so we aren't exactly comparing like with like.

Try running this again. Would you expect to see exactly the same results? Compare the new results to the old results on the map.

That's all!