

Geo-location Clustering using the k-means Algorithm

Motivation:

Geolocation is the identification or estimation of the real-world geographic location of an object, such as a radar source, mobile phone, or Internet-connected computer terminal ([Wikipedia](#)). Today's world is a digital world and every person is connected with some kind of devices. With geolocation, we can easily identify the location of a person or device. This can be useful in many ways, government and police uses geolocation to track people. Some companies use it for implementing augmented reality in games such as Pokémon go. If a company wants to know the user or website visitor location, it can use geolocation of that user.

Clustering creates a group of data with common features. Geo-location clustering can be useful for analyzing the user data of the company based on user locations and their buying patterns. Based on the geo-location clusters, a company can decide marketing strategies, warehouse locations etc., for individual clusters.

Approach:

1) Getting data and Analyzing data:

There are 3 data files given for this task – DeviceStatus, lat longs, Sample_geo.

- i) DeviceStatus file contains the information collected from devices of a private network. It has features such as device ID, current status, location etc.
- ii) Lat longs file has the data of latitude and longitude co-ordinates from the entire world.
- iii) Sample_geo file contains the subset of latitude-longitude co-ordinates

2) Using AWS service – S3 bucket for storing the data files:

Amazon simple storage service is an object storage service which can accommodate large datasets. It provides scalability, data availability, security and performance.

3) Using AWS service – EMR (Elastic MapReduce) for processing the data:

Amazon Elastic MapReduce is a web services tool for big data processing and analysis. It supports workloads based on Apache Spark, Presto, Hive, Pig functionality.

- 4) Preprocessing the data using Pyspark techniques:
 - i) Extraction: Data files are uploaded to S3 bucket. I extracted the data files from s3 bucket.
 - ii) Transformation: the file extracted from S3 bucket is a RDD object, so in this step I transformed RDD object to spark data frame with appropriate delimiters for each file.
 - iii) Selection: As per the requirement of a task, I have selected few features from file. Also for one file, instead of taking the entire world data, I have used only USA data.
- 5) Creating the clusters for each data file:

Using K-means clustering algorithm, I have tried to cluster the data of latitude and longitude.
- 6) Visualizing the clusters using Basemap:

Basemap is a python library which is mostly used for visualizing latitude-longitude for particular region or entire world. I have used same library for visualizing my analysis.

System Configuration:

Below are the steps I performed as a part of system configuration

- 1) **Create EC2 instance:**
 - (1) Login to AWS educate account and open EC2 dashboard
 - (2) Click on 'Launch Instance' and select 'Ubuntu 18.04' AMI
 - (3) Click 'Review and Launch' button
 - (4) When pop up appears- click on create new key-pair file and give name of key-pair file and download it.
 - (5) Launch the instance
 - (6) For this assignment I have used port# 8888 so in security groups of EC2 instance, I have added this port
- 2) **Create S3 bucket:**
 - (1) Open AWS S3 service console
 - (2) Click on create bucket link
 - (3) Give the name of bucket and hit 'Create' button
- 3) **Setup and configure EMR cluster:**
 - (1) Open EMR console and click on create cluster option
 - (2) Give the name of cluster

- (3) In Application section, select 'Spark' instead of 'Hadoop'
- (4) Select 'm4.xlarge' for instance type
- (5) Give EC2 key-pair name which was created earlier and click create option
- (6) In EC2 instance dashboard, when Hadoop master and slave instances are in running state, edit Security rules for mater EC2 instance to add port# 8888

4) Opening a pyspark jupyter notebook in a browser:

- 1) SSH into a Hadoop master node
- 2) Install all required packages using pip (I have installed boto3, ipyparallel, jupyter, basemap, pyyaml, pandas and ipython)
- 3) Set the pyspark driver path in '.bashrc' file
- 4) In terminal, give command 'pyspark'
- 5) Open a notebook using IPV4 address of master node and token generated in terminal

Big data application/dataset:

1) Description of data processing

- (1) For this task 3 files are given. I have uploaded all 3 files to S3 bucket. Then using boto3, I am reading each file from bucket to spark dataframe.

```
In [8]: s3 = boto3.client('s3',aws_access_key_id='ASIA6DHH57K7XGORCC4N',aws_secret_access_key='sKcdUKdqMNQVpnjKGPnEmR66Fps73hWCm92fj2Ig'
bucket = 'finalprojectdsde'
key = 'devicestatus.txt'

obj = s3.get_object(Bucket=bucket, Key=key)
#df = pd.read_csv(obj['Body'], sep=r'\\|\\|\\|', header=None, engine='python')

In [9]: spark = SparkSession \
.builder \
.appName("Python Spark SQL basic example") \
.config("spark.some.config.option", "some-value") \
.getOrCreate()

In [5]: print(sqlContext)

<pyspark.sql.context.SQLContext object at 0x7f7ab70c2490>

In [10]: body = obj['Body']
csv_file = body.read().decode('utf-8')

#df = spark.read.csv(StringIO(csv_file),header=None, sep=r'\\|\\|\\|')

In [11]: df = pd.read_csv(StringIO(csv_file),header=None, sep=r'\\|\\|\\|', engine='python')

In [12]: df1 = spark.createDataFrame(df)
```

- (2) After reading file into spark dataframe, I have dropped and renamed columns. And also filtered out the data with latitude and longitude values as 0.0

- (3) Then I have performed few more preprocessing steps like rearranging dataframe columns, splitting the columns etc.
- (4) As a final step of data processing, I have uploaded the processed dataframe as a CSV file back to the S3 bucket.

```
In [20]: pandas_df = df3.select("*").toPandas()
```

```
In [19]: file_name = "MobileNetData_spark.csv"
df3.to_csv(file_name)
s3.upload_file(file_name, 'finalprojectdsde', 'devicedata1')
```



```
DeviceStatus - Notepad
File Edit Format View Help
,Latitude,Longitude,Date,DeviceId,Manufacturer,Model
0,33.6894754264,-117.54330825299999,2014-03-15:10:10:20,8cc3b47e-bd01-4482-b500-28f2342679af,Sorrento,F41L
1,37.4321088904,-121.485029632,2014-03-15:10:10:20,ef8c7564-0a1a-4650-a655-c8bbd5f8f943,MeeToo,1.0
2,39.43789083489999,-120.93897848600001,2014-03-15:10:10:20,23eba027-b95a-4729-9a4b-a3cca51c5548,MeeToo,1.0
3,39.363518676700004,-119.4003470799999,2014-03-15:10:10:20,707daba1-5640-4d60-6d9-1d6fa0645be0,Sorrento,F41L
4,33.1913581092,-116.44824264299999,2014-03-15:10:10:20,d666fe81-aa55-43b4-9418-fc6e7a00f891,Ronin,Novelty Note 1
5,33.8343543748,-117.330008857,2014-03-15:10:10:20,ffa18088-69a0-433e-84b8-006b2b9cc1d0,Sorrento,F41L
6,37.3803954321,-121.840756755,2014-03-15:10:10:20,66d678e6-9c87-48d2-a415-8d5035e54a23,Sorrento,F33L
7,34.1841062345,-117.94353290000001,2014-03-15:10:10:20,673f7e4b-d52b-44fc-8826-aea460c3481a,MeeToo,4.1
8,32.2850556785,-111.81958373399999,2014-03-15:10:10:20,a678ccc3-b0d2-452d-bf89-85bd095e28ee,Ronin,Novelty Note 2
9,45.2400522984,-122.377467861,2014-03-15:10:10:20,86bef6ae-2f1c-42ec-aa67-6acced7b0675,Sorrento,F41L
10,37.9248961741,-122.20686816700001,2014-03-15:10:10:20,27178d24-3a61-42f7-a784-e3263f25cc6f,iFruit,3
11,38.1653163975,-122.1516083799999,2014-03-15:10:10:20,b4a15931-9a69-469f-9823-a45974472c51,Titanic,2400
12,33.323126641,-116.472234745,2014-03-15:10:10:20,e75dc777-b531-4dbd-88d5-39c77266e6ea,Ronin,S1
13,33.1774985363,-116.889226299,2014-03-15:10:10:20,d4ebd9ae-4dad-4fb4-ba1d-d2a9610a458d,Ronin,Novelty Note 3
14,32.2083493316,-111.43410271299999,2014-03-15:10:10:20,b954db08-1f97-4311-8d42-1a7ba39d8dcf,Ronin,Novelty Note 1
15,34.0487620041,-111.928871717,2014-03-15:10:10:20,16085fbf-cda5-4489-84b9-0fad888f9e7a,MeeToo,1.0
16,37.9031053656,-121.56145134200001,2014-03-15:10:10:20,6474caf1-7bbf-4594-a526-9ba8ea82e151,iFruit,1
17,36.032967794,-118.970108886,2014-03-15:10:10:20,668e6f06-a8aa-4be5-8609-899c45d3caa8,MeeToo,5.0
18,45.0400810371,-117.858004521,2014-03-15:10:10:20,6d195272-8dba-42d6-9b1f-fe61edf7f2a2,Ronin,Novelty Note 1
19,35.2338863976,-114.30575230000001,2014-03-15:10:10:20,d228cdab-8b35-4733-9f2e-e4760dfac3b0,Sorrento,F10L
20,32.820298768,-110.862165369,2014-03-15:10:10:20,92f0ee6e-a56a-4e3c-808f-a6282c539ef8,Sorrento,F41L
21,34.218217050199996,-118.121882092,2014-03-15:10:10:20,89f630b4-7d95-4f3c-aa83-d203b874ade6,Sorrento,F41L
22,38.0180756036,-120.067860245,2014-03-15:10:10:20,fe4674a7-0632-494a-aa3f-f5865a724e1c,iFruit,3A
23,34.1259993213,-117.91415019799999,2014-03-15:10:10:20,ab508fbf-f5ca-4ee6-92f7-c5e090752d2d,Sorrento,F01L
24,0.0,0.0,2014-03-15:10:10:20,90d30931-6664-4a1f-966c-ab0c3b37156e,Ronin,S1
25,32.8247811394,-116.870394352,2014-03-15:10:10:20,9c82628f-cb22-4ceb-998f-e943d1ca5bb,Ronin,S1
26,45.326414381999996,-117.8078110299999,2014-03-15:10:10:20,08bf61ec-f224-4e8c-a754-1ed381329ed4,Titanic,2000
27,39.4708861702,-119.65992609700001,2014-03-15:10:10:20,86e4bd60-5c72-4558-a019-caf1de9d14f1,iFruit,5
28,39.5370541885,-114.751908093,2014-03-15:10:10:20,e205ee4a-e7f8-44ba-bfb0-2f4a7a604e09,MeeToo,1.0
29,39.0791659375,-119.529051299,2014-03-15:10:10:20,6861961e-50e5-4e67-ab4e-36cc327ec6e2,Sorrento,F23L
30,0.0,0.0,2014-03-15:10:10:20,7017f788-0de7-4218-a1bb-1b287944b34,Ronin,Novelty Note 1
31,33.851446595300004,-117.787423338,2014-03-15:10:10:20,3c6ad768-a510-43b4-9050-aad293b75373,Sorrento,F21L
32,33.4467594034,-111.36532692200001,2014-03-15:10:10:20,8316b507-7620-47aa-b56b-cae5cb2cd819,MeeToo,3.0
33,33.7276654194,-112.04447537700001,2014-03-15:10:10:20,5d08d989-c873-4598-a3d1-19cf8f0984fd,Sorrento,F41L
34,34.5711112129,-117.805608832,2014-03-15:10:10:20,fce6f168-7cca-4c00-97f5-9d03ad8232a9,Sorrento,F41L
```

2) Implementation of k-means algorithm:

- (1) For this task, I am reading preprocessed file from S3 bucket in pyspark
- (2) For calculating distance of point from centroids, I have used Euclidean distance and Great Circle distance
- (3) The input parameters – distance function and number of clusters are taken from Command line

```

take_input = raw_input("Enter The Distance function to use and number of clusters : \n")
take_input = take_input.split(" ")

if take_input[0] == "eu":
    dist_funct = EuclideanDistance
elif take_input[0] == "gc":
    dist_funct = GreatCircleDistance
else:
    print("wrong function shoule be 'eu' or 'gc'")
    exit(-1)

centers, clusters = kmeans(filtered_data, dist_funct, int(take_input[1]))

```

- (4) As a part of k-means cluster algorithm, I am arbitrarily choosing k centers. Then using method given as input from CLI, I am calculating distance of each point with center and creating clusters
- (5) In next step, I am recursively updating center points and clusters. I am repeating this process until the difference between old center point and new center point is not as small as 0.01 (convergence criterion)

```

In [34]: def kmeans(filtered_data, dist_funct, k):
    clusters = []
    centers = filtered_data.takeSample(False, k)
    diff = 1.0
    conv_criteria = 0.01
    kMean_df = pd.DataFrame(columns = ['Latitude', 'Longitude', 'cluster_id' ])

    while diff > conv_criteria:
        cluster = filtered_data.map(lambda p: (closestPoint(p, centers, dist_funct), [p, 1]))
        newCenters = cluster.reduceByKey(lambda p1,p2: (addPoints(p1[0],p2[0]), p1[1]+p2[1])).map(lambda l: (l[0], np.array(l[1]
        diff = getDiff(newCenters, centers)
        for (i, newCt) in newCenters:
            centers[i] = newCt

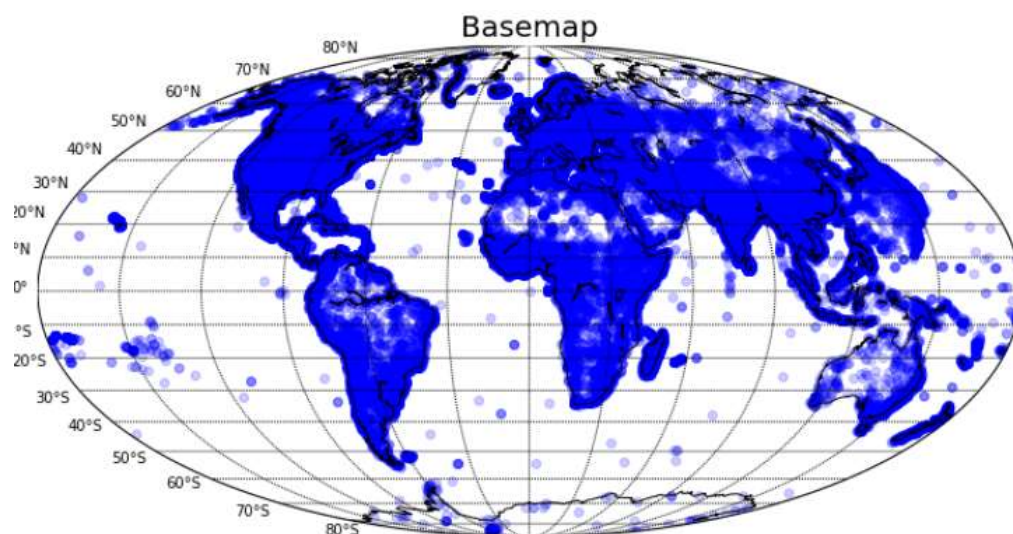
    print(centers)
    return centers, kMean_df

```

3) Visualizing Result of k-means:

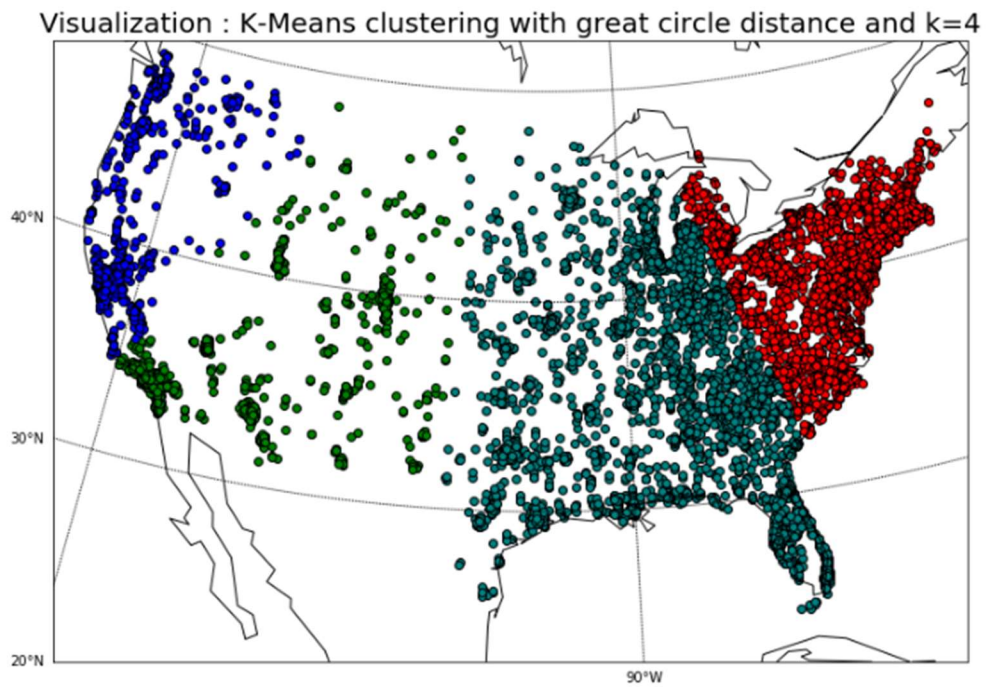
For visualizing all the graphs, I have used basemap library.

- 1) Visualization of latitude-longitude data before implementing k-means clusters:

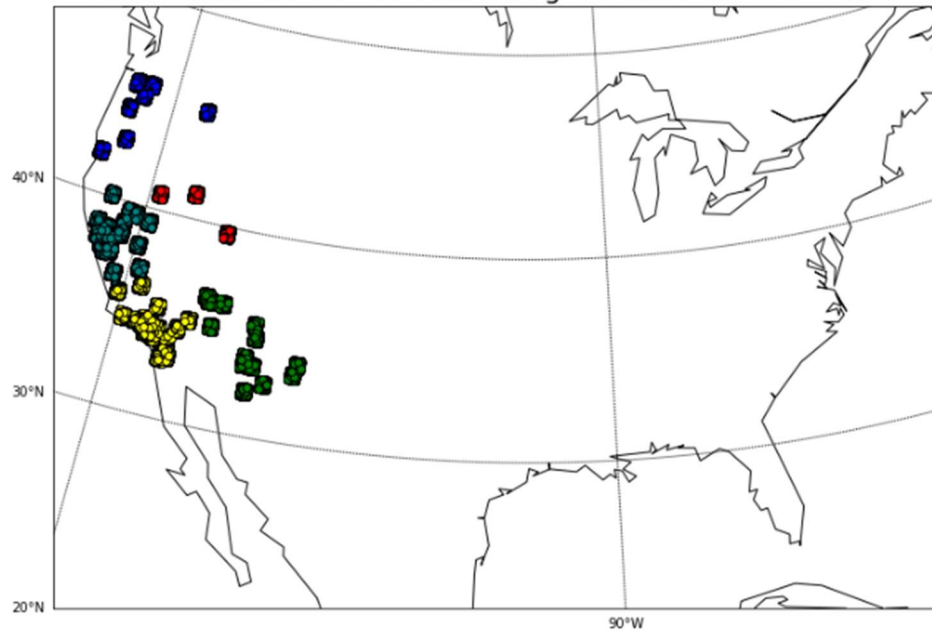


This plot is plotted using all the data of entire world latitude-longitude coordinates.

2) Visualization of data after implementing k-means clusters:



K-means cluster for Device Data with great circle distance and $k = 5$



Conclusion or Future Work:

The latitude-longitude dataset was big enough but with AWS spark, it was easy to do the processing. The process of creating clusters using parallelism was less time consuming.

I have done clustering only for US data. But it can be done on even larger data and data of multiple companies instead of only one company.