

Geo-Location Clustering Using K-Means Algorithm

Introduction:

Clustering, as defined by Stanford University, is a process of grouping a set of data points into clusters such that points within the same cluster are positioned comparable to each other while points from various clusters are different.

Clustering has many useful applications such as finding a group of consumers with common preferences, grouping documents based on the similarity of their contents, or finding spatial clusters of customers to improve logistics. More specific use cases are

- Marketing: given a large set of customer transactions, find customers with similar purchasing behaviors .
- Document classification: cluster web log data to discover groups of similar access patterns.
- Logistics: find the best locations for warehouses or shipping centers to minimize shipping times.

We will approach the clustering problem by implementing the k-means algorithm. k-means is a distance-based method that iteratively updates the location of k cluster centroids until convergence. The main user-defined ingredients of the k-means algorithm are the distance function (often Euclidean distance) and the number of clusters k. This parameter needs to be set according to the application or problem domain. (There is no magic formula to set k.) In a nutshell, k-means groups the data by minimizing the sum of squared distances between the data points and their respective closest centroid.

System Configuration:

S3 Bucket Creation:

Amazon S3

Buckets (3)

Buckets are containers for data stored in S3. [Learn more](#)

Find buckets by name

	Name	Region	Access
<input type="radio"/>	aws-logs-470999077021-us-east-1	US East (N. Virginia) us-east-1	Objects can be public
<input type="radio"/>	device1status	US East (N. Virginia) us-east-1	Objects can be public
<input type="radio"/>	geoclusterresults	US East (N. Virginia) us-east-1	Objects can be public

Creating EMR Cluster:

In the process of creating EMR cluster, created a key-pair, which then downloaded as .pem file. Using Putty SSH client created .ppk file from .pem file.

Key pairs (1)

Filter key pairs

	Name	Fingerprint	ID
<input type="checkbox"/>	geo cluster	c7:f1:f3:ba:49:be:ee:c4:24:11:7b:96:0b...	key-093db56ca3436b2ca

aws Services

Amazon EMR

Clusters

Notebooks

Git repositories

Security configurations

Block public access

VPC subnets

Events

Help

What's new

Clone Terminate AWS CLI export

Cluster: My cluster **Waiting** Cluster ready after last step completed.

Summary Application user interfaces Monitoring Hardware Configurations Events Steps Bootstrap actions

Summary

ID: j-361YUTM9RZVRL

Creation date: 2020-12-06 08:43 (UTC-5)

Elapsed time: 2 days, 8 hours

After last step completes: Cluster waits

Termination protection: Off [Change](#)

Tags: -- [View All / Edit](#)

Master public DNS: ec2-54-209-5-184.compute-1.amazonaws.com [Connect to the Master Node Using SSH](#)

Configuration details

Release label: emr-5.32.0

Hadoop distribution: Amazon

Applications: Spark 2.4.7, Zeppelin 0.8.2

Log URI: s3://aws-logs-470999077021-us-east-1/elasticmapreduce/ [View Log](#)

EMRFS consistent view: Disabled

Custom AMI ID: --

Application user interfaces

Persistent user interfaces [View](#): [Spark history server](#), [YARN timeline server](#)

On-cluster user interfaces [View](#): Not Enabled [Enable an SSH Connection](#)

Network and hardware

Availability zone: us-east-1e

Subnet ID: [subnet-2c4e9b1d](#)

Master: **Running** 1 m4.xlarge

Core: **Running** 2 m4.xlarge

Task: --

Cluster scaling: Not enabled

Security and access

Key name: geo cluster

EC2 instance profile: EMR_EC2_DefaultRole

EMR role: EMR_DefaultRole

Data Preparation :

Data Preparation Before implementing the actual algorithm, we went through pre-processing step in order to convert the data into a standardized format for later processing. The following describes the pre-processing process for device status data:

1. Load the dataset
2. Determine which delimiter to use
3. Filter out any records which do not parse correctly; each record should have exactly 14 values
4. Extract the date, model, device ID, and latitude and longitude a. date: 1st field b. model: 2nd field c. device ID: 3rd field d. latitude: 13th field e. longitude: 14th field.
5. Store latitude and longitude as the first two fields
6. Filter out locations that have a latitude and longitude of 0
7. Split the model field that contains the device manufacturer and model name by spaces
8. Saved the extracted data to comma delimited text files in S3.

Location: 's3://geoclusterresults/devicedata.csv'

Data:

Mobilenet:

```
In [23]: #Printing the data frame after dropping zeros
devicedata
```

Out[23]:

	latitude	longitude	date	model	device ID
0	33.6894754264	-117.543308253	2014-03-15:10:10:20	F41L	8cc3b47e-bd01-4482-b500-28f2342679af
1	39.3635186767	-119.400334708	2014-03-15:10:10:20	F41L	707daba1-5640-4d60-a6d9-1d6fa0645be0
2	33.1913581092	-116.448242643	2014-03-15:10:10:20	Novelty Note 1	db66fe81-aa55-43b4-9418-fc6e7a00f891
3	33.8343543748	-117.330000857	2014-03-15:10:10:20	F41L	ffa18088-69a0-433e-84b8-006b2b9cc1d0
4	37.3803954321	-121.840756755	2014-03-15:10:10:20	F33L	66d678e6-9c87-48d2-a415-8d5035e54a23
...
65640	39.4463417571	-114.736213453	2014-03-15:10:49:30	F22L	40e61459-5448-4dc9-bb89-42e73a4e19cf
65641	38.4282665514	-121.25933863	2014-03-15:10:49:30	S2	b13ece99-62ab-4c9f-a366-6a06bd5e877f
65642	33.7778202246	-108.575470704	2014-03-15:10:49:30	F41L	32af1a0b-ca7f-4906-9772-9eb9435e7e4c
65643	38.2596913494	-122.295712621	2014-03-15:10:49:30	S1	a48a5559-d916-481b-84a9-5dce6272cce1
65644	34.2415255221	-118.23526739	2014-03-15:10:49:30	2	d86fbaa6-b71b-435f-a0bf-5304a202a70b

65645 rows × 5 columns

Synthetic:

```
In [38]: #Displaying Sample_geo data
geo_2.head(5)
```

Out[38]: [Row(Latitude='37.77253945', Longitude='-77.49954987', LocationID='1'),
Row(Latitude='42.09013298', Longitude='-87.68915558', LocationID='2'),
Row(Latitude='39.56341754', Longitude='-75.58753204', LocationID='3'),
Row(Latitude='39.45302347', Longitude='-87.69374084', LocationID='4'),
Row(Latitude='38.9537989', Longitude='-77.01656342', LocationID='5')]

DBpedia:

:

	lat		long	name_of_page
0	36.7	3.216666666666667		<http://dbpedia.org/resource/Algeria>
1	42.5	1.5166666666666666		<http://dbpedia.org/resource/Andorra>
2	12.516666666666667	-70.03333333333333		<http://dbpedia.org/resource/Aruba>
3	-8.833333333333334	13.333333333333334		<http://dbpedia.org/resource/Angola>
4	41.333333333333336		19.8	<http://dbpedia.org/resource/Albania>

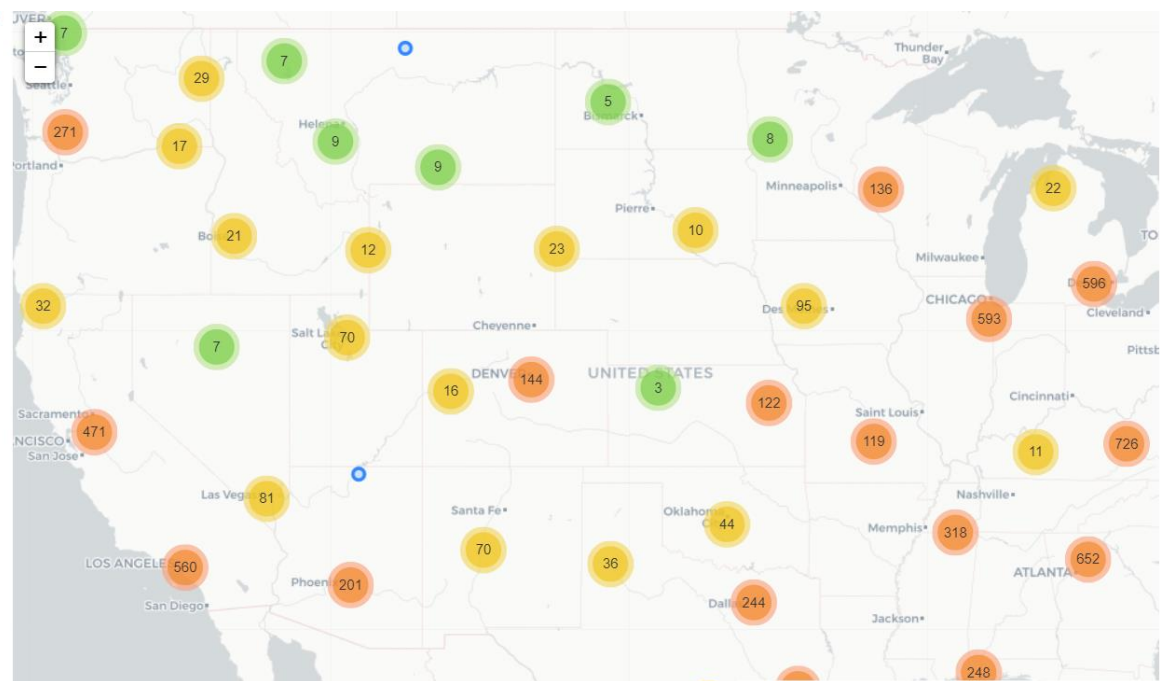
Data Visualization:

Mobilenet:

Map showing the distribution of the number of people with the last name 'Smith' by county in the Western United States. The map includes California, Nevada, Utah, Arizona, and parts of Oregon, Idaho, and New Mexico. The size of the orange circle indicates the number of people with the last name 'Smith' in that county.

County	Number of people with last name 'Smith'
Alameda	2361
Alameda	555
Alameda	1051
Alameda	10522
Alameda	4686
Alameda	2760
Alameda	482
Alameda	724
Alameda	995
Alameda	1914
Alameda	2727
Alameda	19564
Alameda	4194
Alameda	987
Alameda	3543
Alameda	1049

Out[44]:



Out[13]:

