

- خوشه‌بندی آفلاین

به‌منظور خوشه‌بندی آفلاین داده‌ها، از بستر PySpark استفاده شده است. با توجه به حجم زیاد داده‌ها، استفاده از این بستر موجب افزایش سرعت عملیات خوشه‌بندی داده‌ها می‌شود. در ادامه گام‌های زیر انجام می‌شود:

❖ پس از تعریف spark session، با استفاده از دستور نشان داده شده در شکل زیر داده‌ها خوانده می‌شوند.

```
data = sc.read \
    .option('header', 'True')\
    .option('inferSchema', 'True')\
    .option('sep', ',')\
    .csv('/content/drive/MyDrive/data.csv')
```

❖ در گام بعدی با استفاده از قطعه کد قابل مشاهده در شکل زیر، حجم ۲۰ درصد از کل داده‌ها محاسبه می‌شود.

```
import math

data_len = data.count()
length = math.trunc(data_len * 0.2)
```

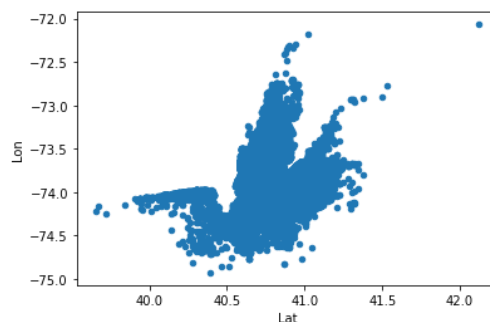
❖ ویژگی‌هایی که برای خوشه‌بندی داده‌ها مورد استفاده قرار گرفته‌اند، ویژگی‌های longitude و latitude می‌باشند. بنابراین سایر ستون‌ها از دیتافریم حذف می‌شوند. سپس با کمک تابع distinct داده‌های تکراری حذف می‌شوند.

```
data = data.na.drop()

# pick Longitude and Latitude features
unimportant_features = ['Date/Time', 'Base']
data = data.drop(*tuple(unimportant_features))

df = data.distinct()
```

❖ در گام بعدی بیست درصد داده‌ها جداسازی شده و شافل می‌شوند. شکل زیر نشان‌دهنده‌ی نمودار پراکندگی داده‌های مورد استفاده برای خوشه‌بندی آفلاین است.



❖ در ادامه با استفاده از vectorizer و standardScaler مقادیر ستون‌ها به بردار تبدیل شده و نرمال‌سازی می‌شوند.

```
vectorize latitude and longitude features

from pyspark.ml.feature import VectorAssembler

assemble=VectorAssembler(inputCols=[
    'Lat', 'Lon'], outputCol='features')

assembled_data=assemble.transform(new_df)
assembled_data.show(2)

standardize data

[ ] from pyspark.ml.feature import StandardScaler
scale=StandardScaler(inputCol='features',outputCol='standardized')
data_scale=scale.fit(assembled_data)
data_scale_output=data_scale.transform(assembled_data)
data_scale_output.show(2)
```

❖ الگوریتم خوشه‌بندی مورد استفاده در این فاز، الگوریتم k-means می‌باشد. به‌منظور یافتن بهترین تعداد خوشه، الگوریتم به‌ازای مقادیر مختلف k اجرا می‌شود. سپس به‌ازای هر خوشه‌بندی مقادیر <sup>1</sup>SSE و silhouette به‌طور مجزا ذخیره می‌شوند. در ادامه الگوریتم خوشه‌بندی به‌ازای دو تا بیست خوشه اجرا شده است.

```
from pyspark.ml.clustering import KMeans
from pyspark.ml.evaluation import ClusteringEvaluator
silhouette_score=[]
SSE = []

evaluator = ClusteringEvaluator(predictionCol='prediction', featuresCol='standardized', \
    metricName='silhouette', distanceMeasure='squaredEuclidean')

for i in range(2,20):

    KMeans_algo=KMeans(featuresCol='standardized', k=i)

    KMeans_fit=KMeans_algo.fit(data_scale_output)

    output=KMeans_fit.transform(data_scale_output)

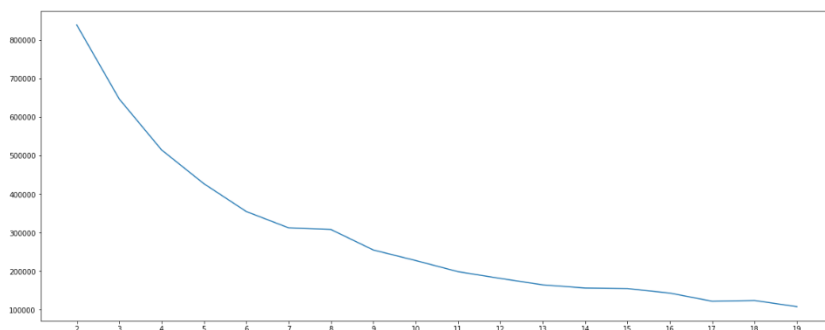
    score=evaluator.evaluate(output)
    silhouette_score.append(score)
    cost = KMeans_fit.summary.trainingCost
    SSE.append(cost)
```

با استفاده از روش آرنج<sup>۲</sup>، می‌توان بهترین مقدار برای متغیر k را تخمین زد. بر مبنای این روش نموداری بر اساس مقادیر خطا به‌ازای هر k رسم می‌شود. نقطه‌ای که نمودار در آن دچار شکست شده باشد، می‌تواند انتخاب مناسبی برای k باشد. شکل زیر نشان‌دهنده‌ی نمودار خطا بر حسب تعداد خوشه است.

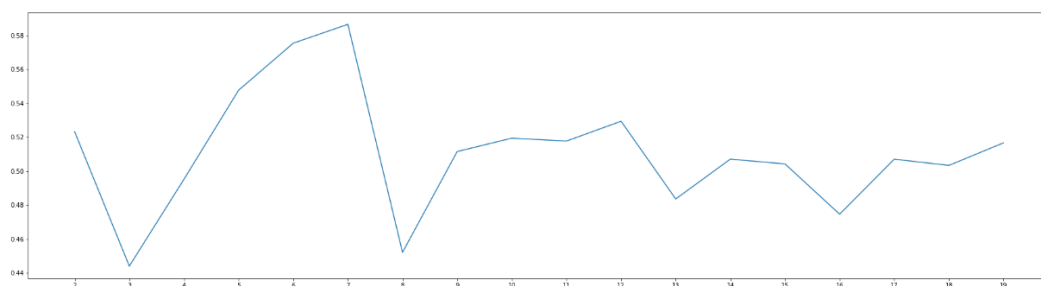
---

<sup>1</sup> Sum of Squared Errors

<sup>2</sup> Elbow Method



با توجه به نمودار بالا مقدار مناسب برای  $k$  یکی از دو مقدار ۷ یا ۸ است زیرا نمودار در این مقدار دچار شکست شده است. به منظور تصمیم‌گیری بهتر در انتخاب تعداد خوشه، نمودار silhouette score بر حسب تعداد خوشه نیز ترسیم شده است. هر قدر معیار silhouette score برای یک خوشه‌بندی بالاتر باشد، آن خوشه‌بندی مناسب‌تر است. با توجه به نمودار به دست آمده، مقدار  $k=7$  بالاترین مقدار silhouette score را دارد. بنابراین عدد هفت به عنوان تعداد خوشه‌ها انتخاب می‌شود.



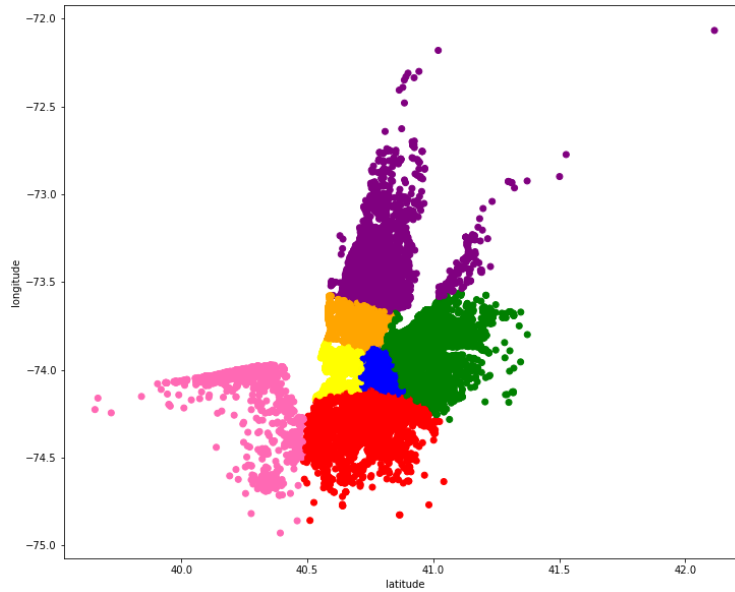
در ادامه خوشه‌بندی به‌ازای هفت خوشه انجام شده و مدل ذخیره می‌شود. با استفاده از قطعه کد زیر می‌توان نمودار نشان‌دهنده‌ی داده‌های خوشه‌بندی شده را نمایش داد.

```
import numpy as np
from pyspark.sql import SQLContext
sqlContext = SQLContext(sc)

transformed = c_pred.select('Lat', 'Lon', 'prediction')
rows = transformed.collect()
df_pred = sqlContext.createDataFrame(rows)
pddf_pred = df_pred.toPandas()

colors=("red","blue","green", "yellow","hotpink","purple","orange","aqua","lightgreen","lightcoral","grey","darkviolet",
        "gold","wheat","rosybrown")

thowdee = plt.figure(figsize=(12,10)).gca()
thowdee.scatter(pddf_pred.Lat, pddf_pred.Lon, c= np.take(colors, pddf_pred.prediction))
thowdee.set_xlabel('latitude')
thowdee.set_ylabel('longitude')
plt.show()
```



## • پرسش داده‌ها در elasticsearch

با استفاده از کوئری‌های زیر، می‌توان موارد خواسته‌شده در صورت پروژه را به‌دست آورد.

```
1 GET /taxi_service_index/_search
2 {
3   "aggs": {
4     "geo": {
5       "geohash_grid": {
6         "field": "Location",
7         "precision": 7,
8         "size": 10
9       },
10      "aggs": {
11        "point": {
12          "geo_centroid": {
13            "field": "Location"
14          }
15        }
16      }
17    }
18  }
19 }
```

```
1 {
2   "took": 546,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 10000,
13      "relation": "gte"
14    },
15    "max_score": 1.0,
16    "hits": [
17      {
18        "_index": "taxi_service_index",
19        "_id": "973",
20        "_score": 1.0,
21        "_source": {
22          "Date/Time": "09/04/2014 10:12:00",
23          "Lat": 40.7356
24        }
25      }
26    ]
27  }
28 }
```

```
1 GET /taxi_service_index/_search
2 {
3   "aggs": {
4     "filter_date": {
5       "filter": {
6         "range": {
7           "Date/Time": {
8             "gte": "09/01/2014 00:01:00",
9             "lte": "09/01/2014 00:20:00"
10          }
11        }
12      },
13      "aggs": {
14        "geo": {
15          "geohash_grid": {
16            "field": "Location",
17            "precision": 7,
18            "size": 10
19          },
20          "aggs": {
21            "point": {
22              "geo_centroid": {
23                "field": "Location"
24              }
25            }
26          }
27        }
28      }
29    }
30  }
```

```
21 {
22   "_source": {
23     "Date/Time": "09/04/2014 10:12:00",
24     "Lat": 40.7356,
25     "Lon": -73.975,
26     "Base": "B02512",
27     "Cluster_number": 0,
28     "Location": {
29       "lat": 40.7356,
30       "lon": -73.975
31     }
32   },
33   {
34     "_index": "taxi_service_index",
35     "_id": "974",
36     "_score": 1.0,
37     "_source": {
38       "Date/Time": "09/04/2014 10:16:00",
39       "Lat": 40.763,
40       "Lon": -73.9699,
41       "Base": "B02512",
42       "Cluster_number": 0,
43       "Location": {
44         "lat": 40.763,
45         "lon": -73.9699
46       }
47     }
48   }
49 }
```

```

1 GET /taxi_service_index/_search
2 {
3   "size": 100,
4   "query": {
5     "match_all": {}
6   },
7   "sort": [
8     {
9       "Date/Time": {
10        "order": "desc"
11      }
12    }
13  ]
14 }

```

```

10 "hits": {
11   "total": {
12     "value": 10000,
13     "relation": "gte"
14   },
15   "max_score": null,
16   "hits": [
17     {
18       "_index": "taxi_service_index",
19       "_id": "8010",
20       "_score": null,
21       "_source": {
22         "Date/Time": "09/09/2014 06:44:00",
23         "Lat": 40.7314,
24         "Lon": -73.9866,
25         "Base": "B02512",

```

## • نوشتن داده‌ها در کانال کافکا (فاز چهارم)

به‌منظور نوشتن داده‌ها در تاپیک مربوط به این فاز (elastic topic) از کد زیر استفاده شده است.

```

def _normal_run(self):
    """
    run method for normal type
    create normal
    :return:
    """
    self._generate_index()
    self._normal_setup()
    data_id = 0

    if self.consumer:
        for each in self.consumer:
            produced_data, key, timestamp = self._produce_answer(each, data_id=data_id)
            data_id += 1
            self._send_data(data=produced_data, key=key, timestamp_ms=timestamp)

    else:
        print("No data in previous phase topic")

```

## • پرسش داده‌ها در کاساندرا

به‌منظور اجرای کوئری‌ها می‌توان از cqlsh در محیط cmd استفاده کرد. اما به‌منظور استفاده از نتایج کوئری‌ها در پایتون لازم است cassandra driver با دستور pip نصب شود. کوئری‌های اجراشده به‌صورت زیر هستند.

```

query = "select * from TaxiServiceKeyspace.midday_table limit 2"
rows = session.execute(query)

```

```

lat = request.form.get("lat")
lon = request.form.get("lon")

point_lat = round(float(lat), 4)
point_lon = round(float(lon), 4)
start_point = (point_lat, point_lon)

session.encoder.mapping[tuple] = session.encoder.cql_encode_tuple
query = "select * from TaxiServiceKeyspace.start_coordinates_table WHERE start=%s"

rows = session.execute(query, (start_point, ))

```

```

lat = request.form.get("lat")
lon = request.form.get("lon")

point_lat = round(float(lat), 4)
point_lon = round(float(lon), 4)

start = request.form.get("start")
end = request.form.get("end")
week_date = (start, end)

session.encoder.mapping[tuple] = session.encoder.cql_encode_tuple
query = "select * from TaxiServiceKeyspace.week_table WHERE week=%s AND lat CONTAINS %s AND lon CONTAINS %s AL
rows = session.execute(query, (week_date, point_lat, point_lon))

```

```

lat_list = rows[0].lat
lon_list = rows[0].lon
index_list = []
flag = False
for i in range(0, len(lat_list)):
    if round(lat_list[i], 4) == point_lat and round(lon_list[i], 4) == point_lon:
        flag = True
        index_list.append(i)

print(flag)
if len(index_list) == 0:
    return "no data found"

else:
    form_data = {}
    form_data['time_data'] = []
    form_data['base'] = []
    for i in index_list:
        form_data['time_data'].append(rows[0].date_time[i])
        form_data['base'].append(rows[0].base[i])

```

## • بازیابی داده‌ها در فلسک

با استفاده از فلسک می‌توان داده‌ها را به‌صورت زیر بازیابی نمود. برای مثال با توجه به شکل زیر کاربر می‌تواند اطلاعات سفرهای انجام‌شده برای یک مختصات خاص را در بازه‌ی زمانی دلخواه به‌دست آورد.

Enter requested information

Latitude
Longitude
Start Date
End Date

Submit

Enter requested information

40.2201

-74.0021

2014-09-01

2014-09-08

Submit

برای مثال در بازه‌ی زمانی مشخص شده توسط کاربر تنها یک سفر با اطلاعات زیر وجود داشته است.

← → 127.0.0.1:5000/second

Visual Table, Online... Logistic Regression... Guide to Using Pre... Comment Classifica... zolfaShefreie/Online... Install Kafka on Win... K Means Clustering... Set Up and Run Ap...

**TIME\_DATA**

['2014-09-01 00:01:00']

**BASE**

['B02512']

همچنین می‌توان سایر کوئری‌ها را نیز به همین صورت اجرا نمود. برای این کار لازم است یک تابع در فلسک ایجاد کرده و آدرس و متد مورد نظر را در آن مشخص نمود.

آدرس گیت‌هاب پروژه:

<https://github.com/zolfaShefreie/Online-Taxi-Service>