

test



default ▼

FINISHED    

```
%pyspark
### Import Libraries
import pandas as pd
from pyspark.sql.types import StringType
from pyspark import SQLContext
import numpy as np # arrays and numerical processing
from sklearn.cluster import KMeans # cluster analysis by partitioning
from sklearn.metrics import silhouette_score as silhouette_score
from __future__ import division, print_function
import plotly.plotly as py
import plotly.graph_objs as go
from datetime import datetime
from pyspark.sql import functions as F
from pyspark.sql.functions import col, udf, unix_timestamp
from pyspark.sql.types import DateType
from pyspark.sql.functions import from_unixtime
from pyspark.sql.functions import col, unix_timestamp, round
```

Took 1 sec. Last updated by anonymous at March 25 2017, 12:55:41 PM.

FINISHED    

```

%spark
road_df1 = sqlContext.read.format("com.databricks.spark.csv").option("header", "true").load
/*.*.csv")
road_df2 = sqlContext.read.format("com.databricks.spark.csv").option("header", "true").load
/*.*.csv", schema=custschema)
road_df3 = sqlContext.read.format("com.databricks.spark.csv").option("header", "true").load
/*.*.csv")

```

Took 10 sec. Last updated by anonymous at March 25 2017, 12:55:58 PM.

FINISHED    

```
%pyspark
road_df.printSchema()
```

```
root
  |-- status: string (nullable = true)
  |-- avgMeasuredTime: string (nullable = true)
  |-- avgSpeed: string (nullable = true)
  |-- extID: string (nullable = true)
  |-- medianMeasuredTime: string (nullable = true)
  |-- TIMESTAMP: string (nullable = true)
  |-- vehicleCount: string (nullable = true)
  |-- _id: string (nullable = true)
  |-- REPORT_ID: string (nullable = true)
```

Took 1 sec. Last updated by anonymous at March 24 2017, 3:32:59 PM.

test




 default ▼



↓

me

Zeppelin

test



```
l-- status: string (nullable = true)
l-- avgMeasuredTime: string (nullable = true)
l-- avgSpeed: string (nullable = true)
l-- extID: string (nullable = true)
l-- medianMeasuredTime: string (nullable = true)
l-- TIMESTAMP: string (nullable = true)
l-- vehicleCount: string (nullable = true)
l-- _id: string (nullable = true)
l-- REPORT_ID: string (nullable = true)
l-- POINT_1_STREET: string (nullable = true)
l-- DURATION_IN_SEC: string (nullable = true)
l-- POINT_1_NAME: string (nullable = true)
l-- POINT_1_CITY: string (nullable = true)
l-- POINT_2_NAME: string (nullable = true)
l-- POINT_2_LNG: string (nullable = true)
l-- POINT_2_STREET: string (nullable = true)
l-- DIST_IN_KMH: string (nullable = true)
```

Took 1 sec. Last updated by anonymous at March 24 2017, 4:07:30 PM.

```
%pyspark
```

FINISHED    

```
road_merged_df = merged_df.select([c for c in merged_df.columns if c not in{'DURATION_IN_SEC'}])
road_merged_df.printSchema()
```

```
#type( rd_df)
```

root

```
l-- status: string (nullable = true)
l-- avgMeasuredTime: string (nullable = true)
l-- avgSpeed: string (nullable = true)
l-- extID: string (nullable = true)
l-- medianMeasuredTime: string (nullable = true)
l-- TIMESTAMP: string (nullable = true)
l-- vehicleCount: string (nullable = true)
l-- _id: string (nullable = true)
l-- REPORT_ID: string (nullable = true)
l-- POINT_1_STREET: string (nullable = true)
l-- POINT_1_NAME: string (nullable = true)
l-- POINT_1_CITY: string (nullable = true)
l-- POINT_2_LNG: string (nullable = true)
l-- POINT_2_STREET: string (nullable = true)
l-- NDT_IN_KMH: string (nullable = true)
l-- POINT_2_POSTAL_CODE: string (nullable = true)
l-- POINT_2_COUNTRY: string (nullable = true)
```

Took 0 sec. Last updated by anonymous at March 24 2017, 4:08:01 PM.

```
%pyspark
```

FINISHED ▶ ↻ 📖 ⚙️

```
rd_df = road_merged_df.select([c for c in road_merged_df.columns if c not in
{'DURATION_IN_SEC', 'DISTANCE_IN_METERS', 'POINT_2_NAME'}])
```

test



FINISHED 

```
%pyspark
```

```
format = "yyyy-MM-dd'T'HH:mm:ss"
rd_df = rd_df.select('avgMeasuredTime', 'avgSpeed', 'TIMESTAMP', 'vehicleCount', 'REPORT_ID', 'POINT_1_LNG',
                    from_unixtime(unix_timestamp('TIMESTAMP', format)).cast("timestamp").alias('d
rd_data_selected = rd_df.select([c for c in rd_df.columns if c not in
{'TIMESTAMP', 'REPORT_ID', 'avgMeasuredTime', 'vehicleCount'}])
rd_data_sel = rd_data_selected.select('avgSpeed', 'POINT_2_LNG', 'POINT_1_LAT', 'POINT_2_LAT', 'P
```

Took 1 sec. Last updated by anonymous at March 24 2017, 6:58:59 PM.

FINISHED ▶ ↗ ↘ ↙ ↚ 📖 ⚙️

```
%pyspark
```

```
#newdf = elevDF.select(year(elevDF.date).alias('dt_year'), month(elevDF.date).alias('dt_month'),
#dayofyear(elevDF.date).alias('dt_dayofy'), hour(elevDF.date).alias('dt_hour'), minute(elevDF.date).alias('dt_minute'),
#second(elevDF.date).alias('dt_second'), weekofyear(elevDF.date).alias('dt_week_no'), unix_timestamp(elevDF.date).alias('dt_int'))
```

Took 0 sec. Last updated by anonymous at March 24 2017, 6:58:39 PM.

FINISHED ▶ ⌵ ⌶ ⚙

```
%pyspark
```

```
rd_data_sel.show()
```

| avgSpeed | POINT_2_LNG | POINT_1_LAT | POINT_2_LAT | POINT_1_LNG | hour |
|----------|-------------------|-------------------|-------------------|--------------------|------|
| 93 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 90 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 84 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 88 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 90 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 88 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 87 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 84 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 88 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 90 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 80 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 79 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 85 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 87 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |
| 88 | 10.10702759027481 | 56.21399075104399 | 56.21740644105506 | 10.145073305557275 | 21 |

Took 4 sec. Last updated by anonymous at March 24 2017, 6:56:35 PM.

FINISHED ▶ 🔍 📖 ⚙️

```
%pyspark
```

```
#rd_data_selected.show()
```

```
temp = rd_data_selected.select('avgSpeed','date',F. weekofyear('date').alias('weekofyear'))
```

default ▼

↓

↓

↓

test



Took 0 sec. Last updated by anonymous at March 23 2017, 7:39:40 AM.

Took 1 min 57 sec. Last updated by anonymous at March 24 2017, 4:13:19 PM.

Took 2 min 17 sec. Last updated by anonymous at March 24 2017, 4:41:41 PM.

Took 1 sec. Last updated by anonymous at March 24 2017, 1:12:25 AM.

Took 2 sec. Last updated by anonymous at March 24 2017, 1:12:44 AM.

Took 2 sec. Last updated by anonymous at March 24 2017, 1:13:34 AM.

↓

Zeppelin

test

FINISHED

default

```
%r
tbl1<-subset(data1, data1$km_group == 1 & data1$mean < 30)
#tbl1<-subset(tbl, data1$mean > 50)
result4<-kmeans(tbl1,4)
tbl1$cluster <- result4$cluster
#View(tbl1)
tbl1$cluster <- as.factor(tbl1$cluster)
mapgilbert <- get_map(location = c(lon = mean(tbl1$long1), lat = mean(tbl1$lat1)), zoom = 12,
# plotting the map with some points on it
ggmap(mapgilbert) + geom_point(data = tbl1, aes(x = as.numeric(long1), y = as.numeric(lat1), color = 20) + guides(fill=FALSE, alpha=FALSE, size=FALSE)
```



Took 5 sec. Last updated by anonymous at March 24 2017, 5:13:30 PM.

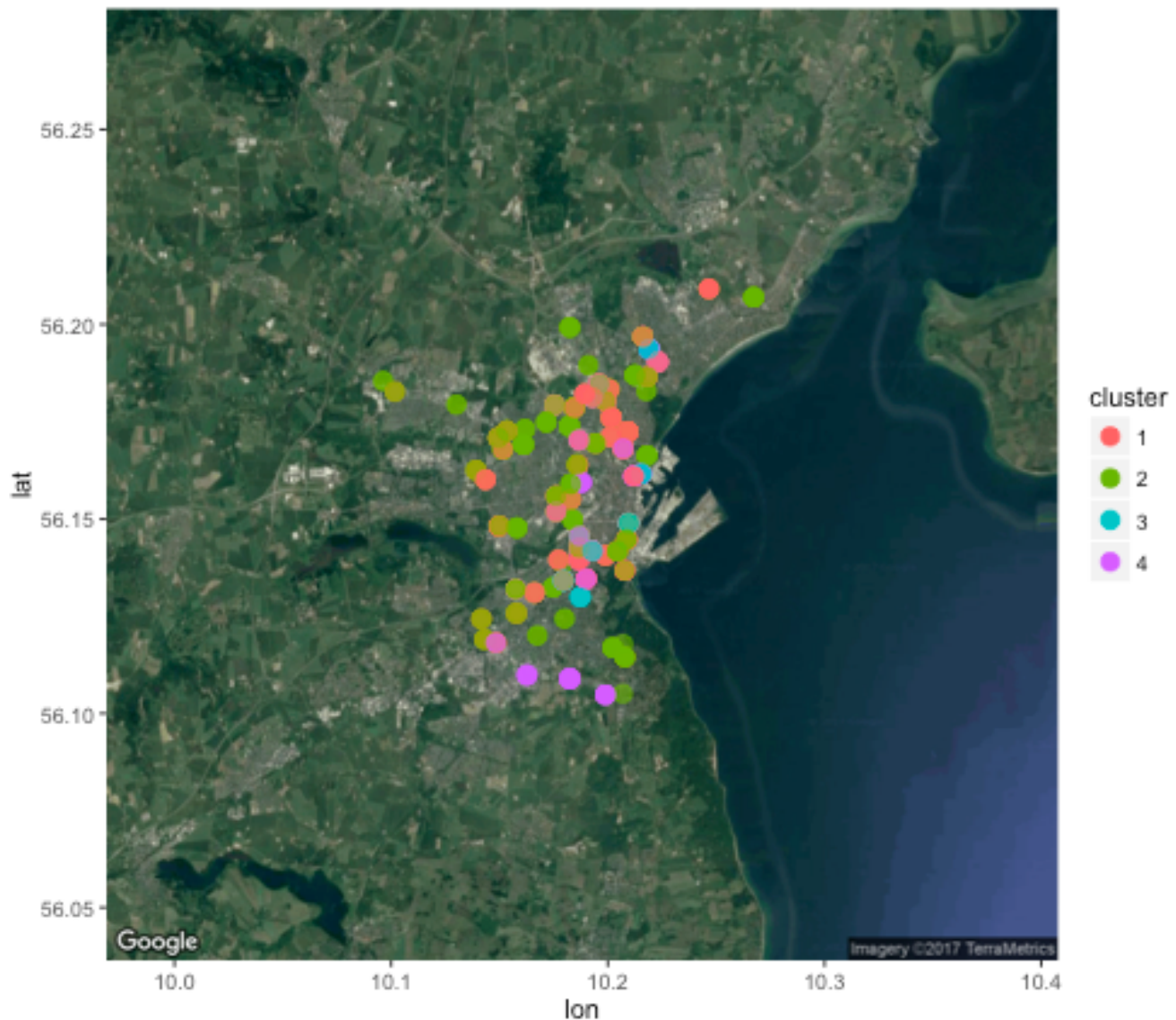
test

```
tbl2<-subset(data1, data1$hr_group == 2 & data1$mean < 30)
result4<-kmeans(tbl2,4)
tbl2$cluster<-result4$cluster
#View(tbl2)
tbl2$cluster <- as.factor(tbl2$cluster)
mapgilbert <- get_map(location = c(lon = mean(tbl2$long1), lat = mean(tbl2$lat1)), zoom = 11,
# plotting the map with some points on it
ggmap(mapgilbert) + geom_point(data = tbl2, aes(x = as.numeric(long1), y = as.numeric(lat1), c
= 20) + guides(fill=FALSE, alpha=FALSE, size=FALSE)
```



Took 6 sec. Last updated by anonymous at March 24 2017, 5:14:06 PM.


```
tbl3$cluster <- result4$cluster
#View(tbl3)
tbl3$cluster <- as.factor(tbl3$cluster)
mapgilbert <- get_map(location = c(lon = mean(tbl3$long1), lat = mean(tbl3$lat1)), zoom = 11,
# plotting the map with some points on it
ggmap(mapgilbert) + geom_point(data = tbl3, aes(x = as.numeric(long1), y = as.numeric(lat1), c
shape = 20) + guides(fill=FALSE, alpha=FALSE, size=FALSE)
```



Took 5 sec. Last updated by anonymous at March 24 2017, 5:14:52 PM.

```
%r
tbl4<-subset(data1, data1$hr_group == 4 & data1$mean < 30)
result4<-kmeans(tbl4,4)
tbl4$cluster <- result4$cluster
```

```
library(ggmap)
library(ggplot2)
library(leaflet)
library(RColorBrewer)

# get the map
mapgilbert <- get_map(location = c(lon = mean(tbl4$long1), lat = mean(tbl4$lat1)), zoom = 12,
# plotting the map with some points on it
ggmap(mapgilbert) + geom_point(data = tbl4, aes(x = as.numeric(long1), y = as.numeric(lat1), color = cluster), size = 20) + guides(fill=FALSE, alpha=FALSE, size=FALSE)
```

Error in check_for_XQuartz(): X11 library is missing: install XQuartz from xquartz.macosforge.org

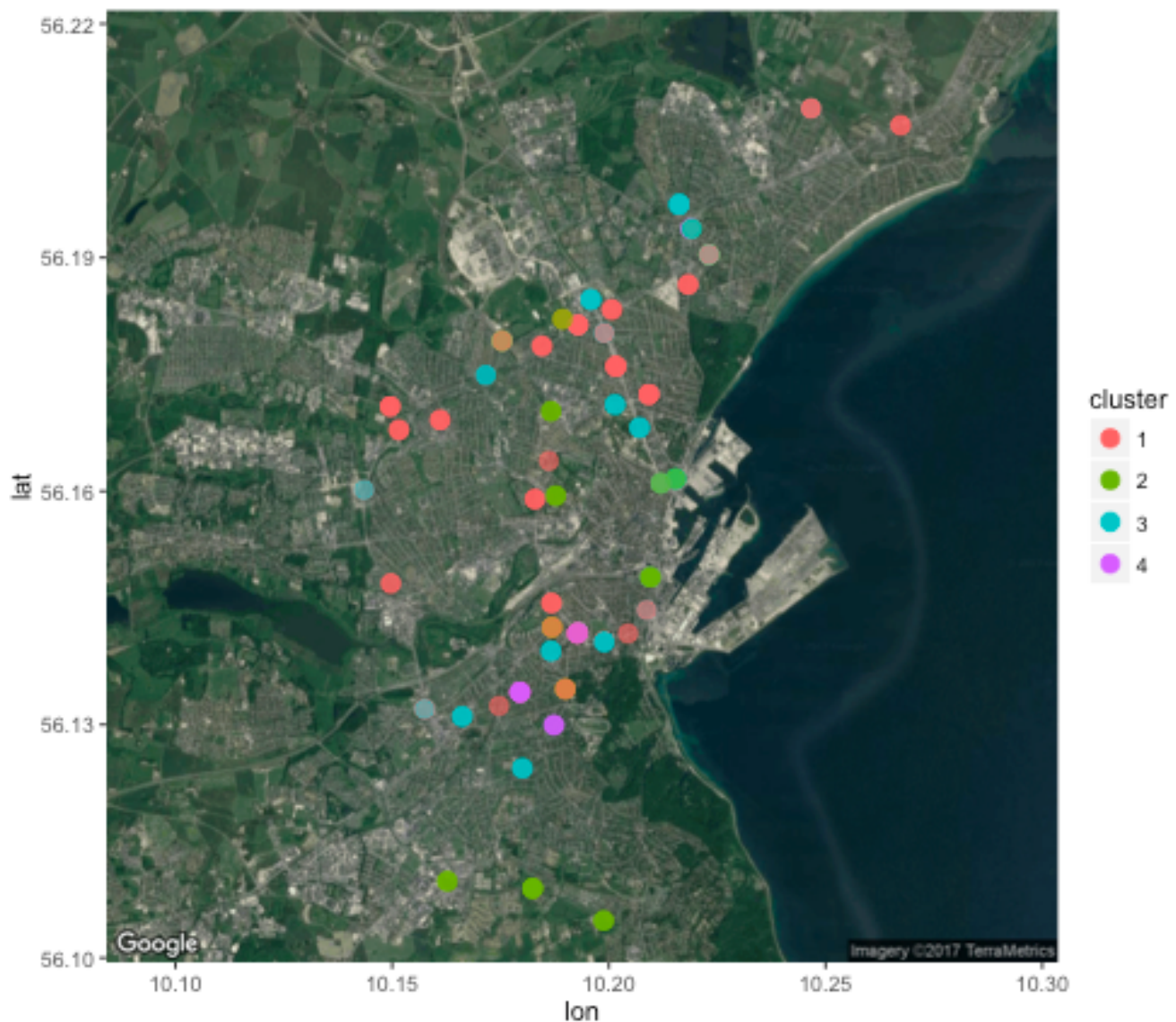


Took 5 sec. Last updated by anonymous at March 24 2017, 5:15:14 PM.

```
%r
tbl5<-subset(data1, data1$hr_group == 5 & data1$mean < 30)
result4<-kmeans(tbl5,4)
```

FINISHED ▶ 🔍 📖 ⚙️

test^{#n1}
game



OK 5 sec. Last updated by anonymous at March 24 2017, 9:19:55 PM.



Zeppelin

READY ▶ ⌵ 📖 ⚙️
↓

test

▶

⌵

📖

✍️

📄

⬇️

📎



default ▾