

NYU CS9223 Big Data Programming

Spring 2017

Assignment 2 – *Yelp Data Challenge*

Report by
Monica Reddy Tirupari
N11349565

1. Summarize the number of reviews by US city, by business category.

Pig Script:

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 A = LOAD './yelp_academic_dataset_business.json' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp: map[]);
4
5 business = FOREACH A GENERATE yelp#'categories' as categories, (int)yelp#'review_count' as review_count, yelp#'city' as city,
6 yelp#'state' as state, (float)yelp#'latitude' as latitude, (float)yelp#'longitude' as longitude;
7
8 coordinates_business = FILTER business BY (latitude<49.384472) AND (latitude>24.520833) AND (longitude<-66.950) AND (longitude>-124.766667);
9
10 us_business = FILTER coordinates_business BY NOT ( (state matches '.*ON.*') OR (state matches '.*QC.*') );
11
12 businesses = FOREACH us_business GENERATE FLATTEN(categories), city, state, review_count;
13
14 grouped = GROUP businesses BY (city,categories);
15
16 result = FOREACH grouped GENERATE group, SUM(businesses.review_count);
17
18 store result into './FinalQ14';

```

Spark Script:

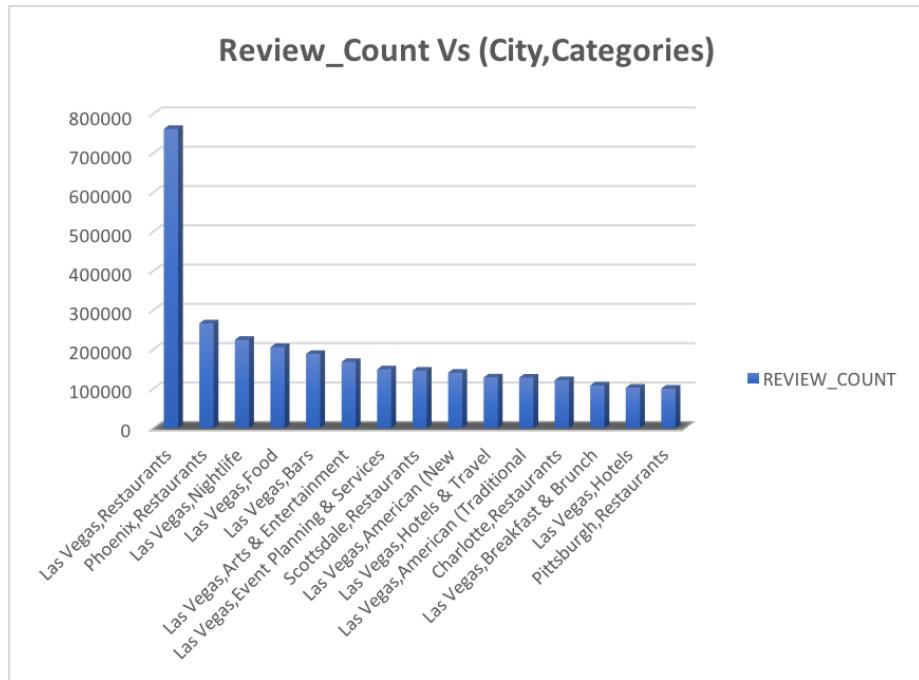
```

Spark1
1 from pyspark.sql import SQLContext
2 sqlContext = SQLContext(sc)
3 df = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_business.json")
4 businesses=df.select(pyspark.sql.functions.explode(
4 df.categories).alias("category"),df.business_id,df.review_count,df.city,df.state,df.longitude,df.latitude)
5 businesses.registerTempTable("business")
6 result=sqlContext.sql("select SUM(review_count) as sum_reviews,city,category from business where latitude >
6 24.520833 and latitude < 49.384472 and longitude > -124.766667 and longitude < -66.950 and state!='ON' and
6 state!='QC' group by city,category")
7 result.write.mode('append').json("/user/cloudera/output_spark/Q1spark.json")
8
9
10

```

Analysis & Visualizations:

	CITY	CATEGORIES	REVIEW_COUNT	CITY,CATEGORIES	REVIEW_COUNT
2	Las Vegas	Restaurants	761661	Las Vegas,Restaurants	761661
3	Phoenix	Restaurants	266766	Phoenix,Restaurants	266766
4	Las Vegas	Nightlife	224955	Las Vegas,Nightlife	224955
5	Las Vegas	Food	206226	Las Vegas,Food	206226
6	Las Vegas	Bars	188813	Las Vegas,Bars	188813
7	Las Vegas	Arts & Entertainment	168321	Las Vegas,Arts & Entertainment	168321
8	Las Vegas	Event Planning & Services	149498	Las Vegas,Event Planning & Services	149498
9	Scottsdale	Restaurants	146271	Scottsdale,Restaurants	146271
10	Las Vegas	American (New)	140884	Las Vegas,American (New)	140884
11	Las Vegas	Hotels & Travel	128908	Las Vegas,Hotels & Travel	128908
12	Las Vegas	American (Traditional)	128502	Las Vegas,American (Traditional)	128502
13	Charlotte	Restaurants	122200	Charlotte,Restaurants	122200
14	Las Vegas	Breakfast & Brunch	108735	Las Vegas,Breakfast & Brunch	108735
15	Las Vegas	Hotels	102928	Las Vegas,Hotels	102928
16	Pittsburgh	Restaurants	100764	Pittsburgh,Restaurants	100764



Assumptions:

- I have considered US border latitudes (24.520833 to 49.384472) and longitudes (-124.766667 to -66.950)
- For a clear visualization, I have only taken Top 15 highest review count tuples of (city, categories) and I have visualized these using excel pivot tables

2. Rank all **cities** by # of stars descending, for **each category**

Pig Script:

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2 A = LOAD './yelp_academic_dataset_business.json' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp: map[]);
3 business = FOREACH A GENERATE yelp#`categories` as categories, (int)yelp#`review_count` as review_count, yelp#`city` as city,
4 yelp#`state` as state, (float)yelp#`latitude` as latitude, (float)yelp#`longitude` as longitude, (float)yelp#`stars` as stars;
5 coordinates_business = FILTER business BY (latitude<49.384472) AND (latitude>24.520833) AND (longitude<-66.950) AND (longitude>-124.766667);
6 us_business = FILTER coordinates_business BY NOT ( (state matches '.*ON.*') OR (state matches '.*QC.*') );
7 businesses = FOREACH us_business GENERATE FLATTEN(categories), city, state, stars;
8 grouped = GROUP businesses BY (categories,city);
9 result = FOREACH grouped GENERATE group, AVG(businesses.stars) as average_stars;
10 Final_result = ORDER result BY group DESC, average_stars DESC;
11 store Final_result into './Q2PIGORDERED';
12
13
14
15
16
17
18
19
20
21
22

```

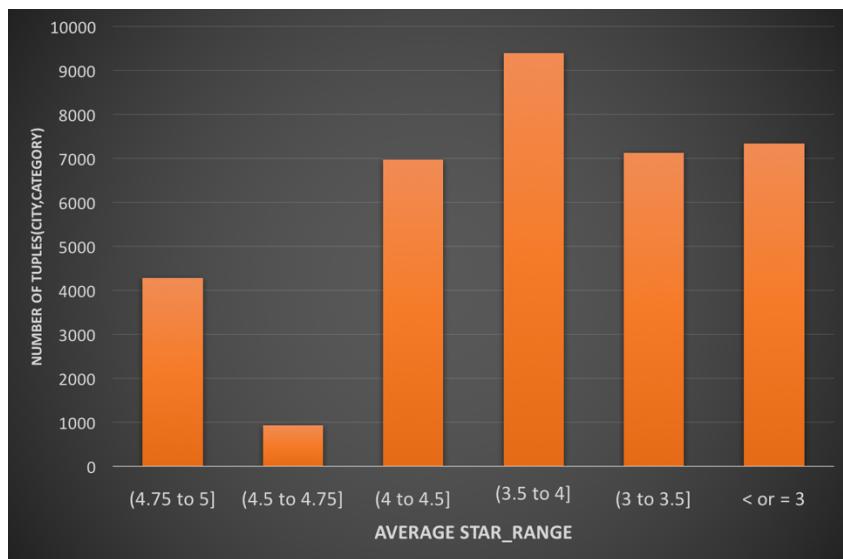
Spark Script:

```

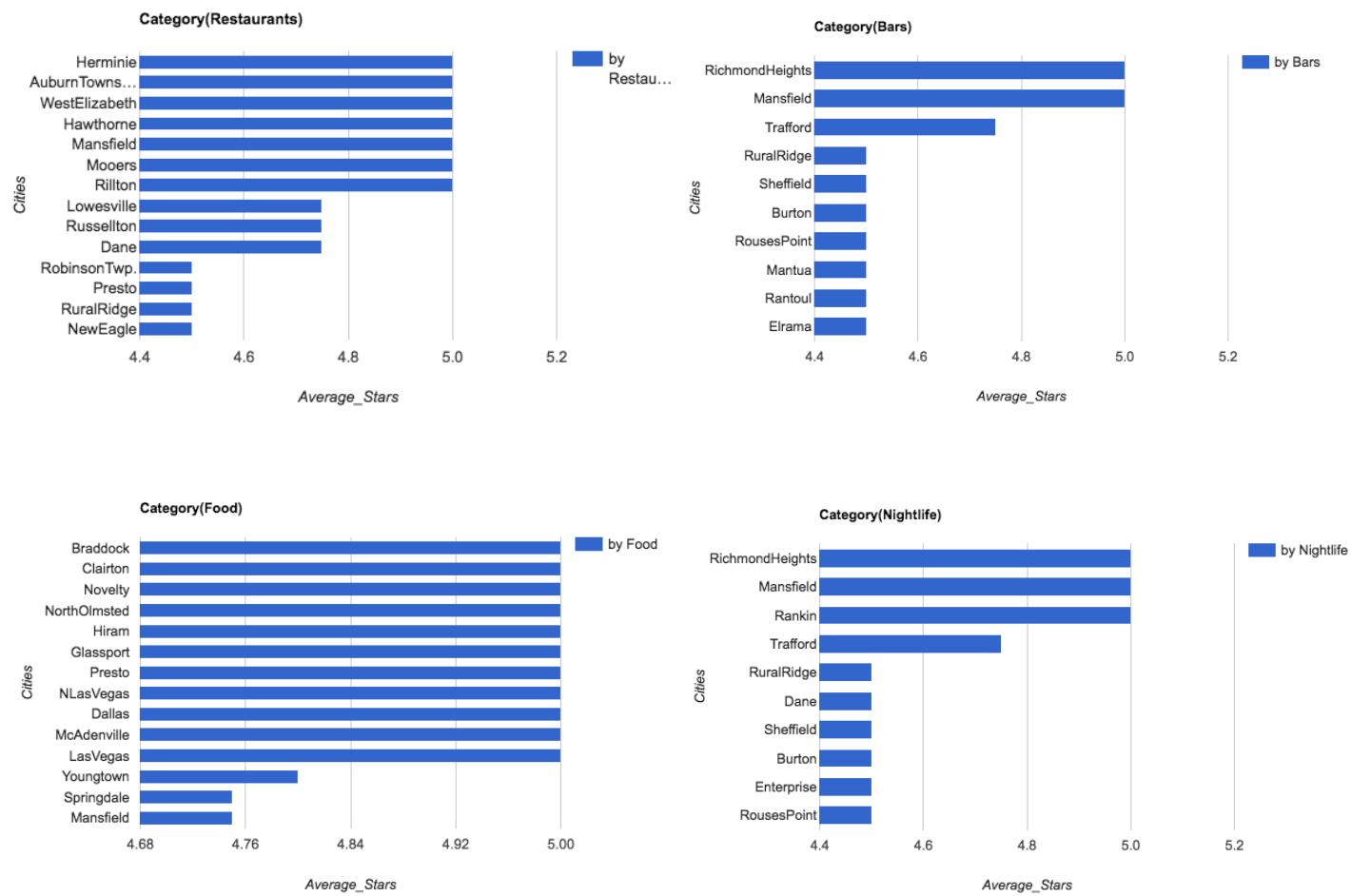
1 from pyspark.sql import SQLContext
2 sqlContext = SQLContext(sc)
3 df = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_business.json")
4 businesses=df.select(pyspark.sql.functions.explode(
5 df.categories).alias("category"),df.business_id,df.stars,df.city,df.state,df.longitude,df.latitude)
6 businesses.registerTempTable("business")
7 result=sqlContext.sql("select AVG(stars) as average_stars,city,category from business where latitude >
8 24.520833 and latitude < 49.384472 and longitude > -124.766667 and longitude < -66.950 and state!= 'ON' and
9 state!= 'QC' group by city,category order by category ASC, average_stars DESC")
10 result.write.mode('append').json("/user/cloudera/output_spark/Q2spark.json")

```

Analysis & Visualizations:



Visualization indicating number of tuples of city, category within an average star range



Assumptions:

- I have considered US border latitudes (24.520833 to 49.384472) and longitudes (-124.766667 to -66.950)

- From the question 1, after analysis i have found that the categories (like restaurants, nightlife, food, bars) have the top most review count. Considering those categories I have visualized average stars for various city having top average star rating.
3. What is the average rank (# stars) for businesses within 10 miles of the University of Wisconsin - Madison, by type of business?

Given Details:

Center: University of Wisconsin - Madison
 Latitude: 43 04' 30" N, Longitude: 89 25' 2" W
 Decimal Degrees: Latitude: 43.0766, Longitude: -89.4125

The bounding box for this problem is ~10 miles, which we will loosely define as 10 minutes. So the bounding box is a square box, 20 minutes long each side (of longitude and latitude), with UWM at the center.

Pig Script:

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 A = LOAD './yelp_academic_dataset_business.json' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp: map[]);
4
5 business = FOREACH A GENERATE yelp#'categories' as categories, (int)yelp#'review_count' as review_count, yelp#'city' as city,
6 yelp#'state' as state, (float)yelp#'latitude' as latitude, (float)yelp#'longitude' as longitude, (float)yelp#'stars' as stars;
7
8 coordinates_business = FILTER business BY (latitude<43.221261) AND (latitude>42.931739) AND (longitude<-89.214450) AND (longitude>-89.610310);
9
10 businesses = FOREACH coordinates_business GENERATE FLATTEN(categories), city, state, stars;
11
12 grouped = GROUP businesses BY (categories, city);
13
14 result = FOREACH grouped GENERATE group, AVG(businesses.stars) as average_stars;
15
16 Final_result = ORDER result BY group DESC, average_stars DESC;
17
18 store Final_result into './FinalQ3PIGORDERED';
19

```

Spark Script:

```

Spark3
1 |from pyspark.sql import SQLContext
2 sqlContext = SQLContext(sc)
3 df = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_business.json")
4 businesses=df.select(pyspark.sql.functions.explode(
5 df.categories).alias("category"),df.business_id,df.stars,df.city,df.longitude,df.latitude)
6 businesses.registerTempTable("business")
7 result=sqlContext.sql("select AVG(stars) as average_stars,city,category from business where latitude > 42.931739 and latitude < 43.221261 and longitude > -89.610310 and longitude < -89.214450 group by city,category order by category ASC, average_stars DESC")
7 result.write.mode('append').json("/user/cloudera/output_spark/Q3spark.json")

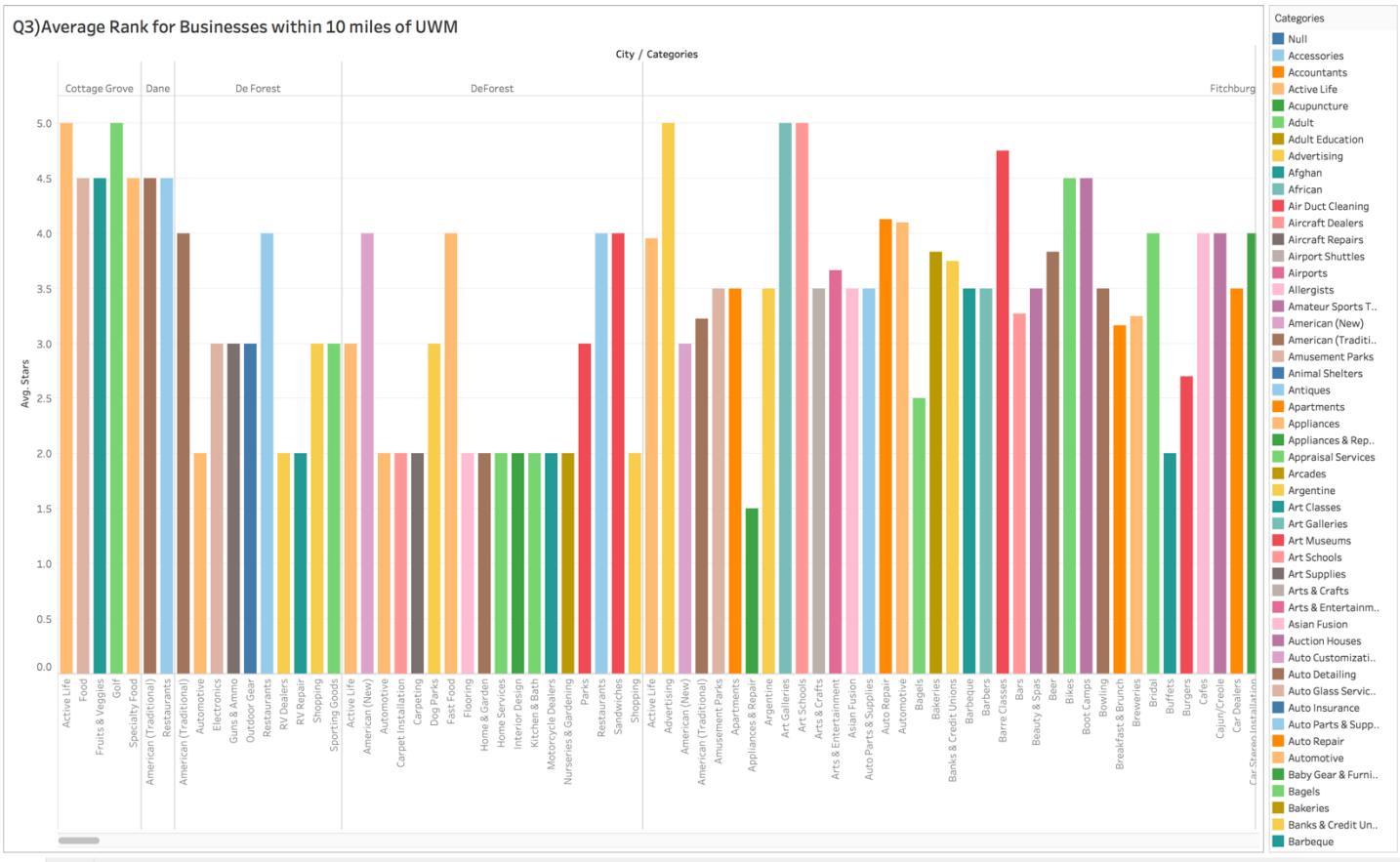
```

Assumptions:

- I have considered 10 miles from center of UWM with latitudes (42.931739 to 43.221261) and longitudes (-89.610310 to -89.214450)

Analysis & Visualizations:

In the Visualization, I have shown Average Star rating for various business categories shown in colors of bar chart within different cities.



4. Rank reviewers by number of reviews. For the top 10 reviewers, show their average number of stars, by category.

Pig Script:

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 A = LOAD './yelp_academic_dataset_business.json' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp: map[]);
4 business = FOREACH A GENERATE yelp#'business_id' as business_id, yelp#'categories' as categories;
5
6 B = LOAD './yelp_academic_dataset_review.json' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp2: map[]);
7
8 review = FOREACH B GENERATE yelp2#'user_id' as user_id1, yelp2#'business_id' as business_id, (int)yelp2#'stars' as stars;
9
10 C = LOAD './yelp_academic_dataset_user.json' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp3: map[]);
11
12 user = FOREACH C GENERATE yelp3#'user_id' as user_id, (int)yelp3#'review_count' as review_count;
13
14 ordered_C = ORDER user BY review_count DESC;
15
16 top_10_C = LIMIT ordered_C 10;
17
18 top_10_stars_joined_BC = JOIN top_10_C BY user_id, review BY user_id1;
19
20 top_10_categories_joined_ABC = JOIN business BY business_id, top_10_stars_joined_BC BY business_id;
21
22 top_10_categories_ABC = FOREACH top_10_categories_joined_ABC GENERATE FLATTEN(categories), stars, user_id;
23
24 grouped = GROUP top_10_categories_ABC BY (user_id, categories);
25
26 grouped_result = FOREACH grouped GENERATE group, AVG(top_10_categories_ABC.stars) as avg_stars;
27
28 ordered_result = ORDER grouped_result BY user_id;
29
30 STORE grouped_result INTO './FinalQ4PIG';
31
32
33 |

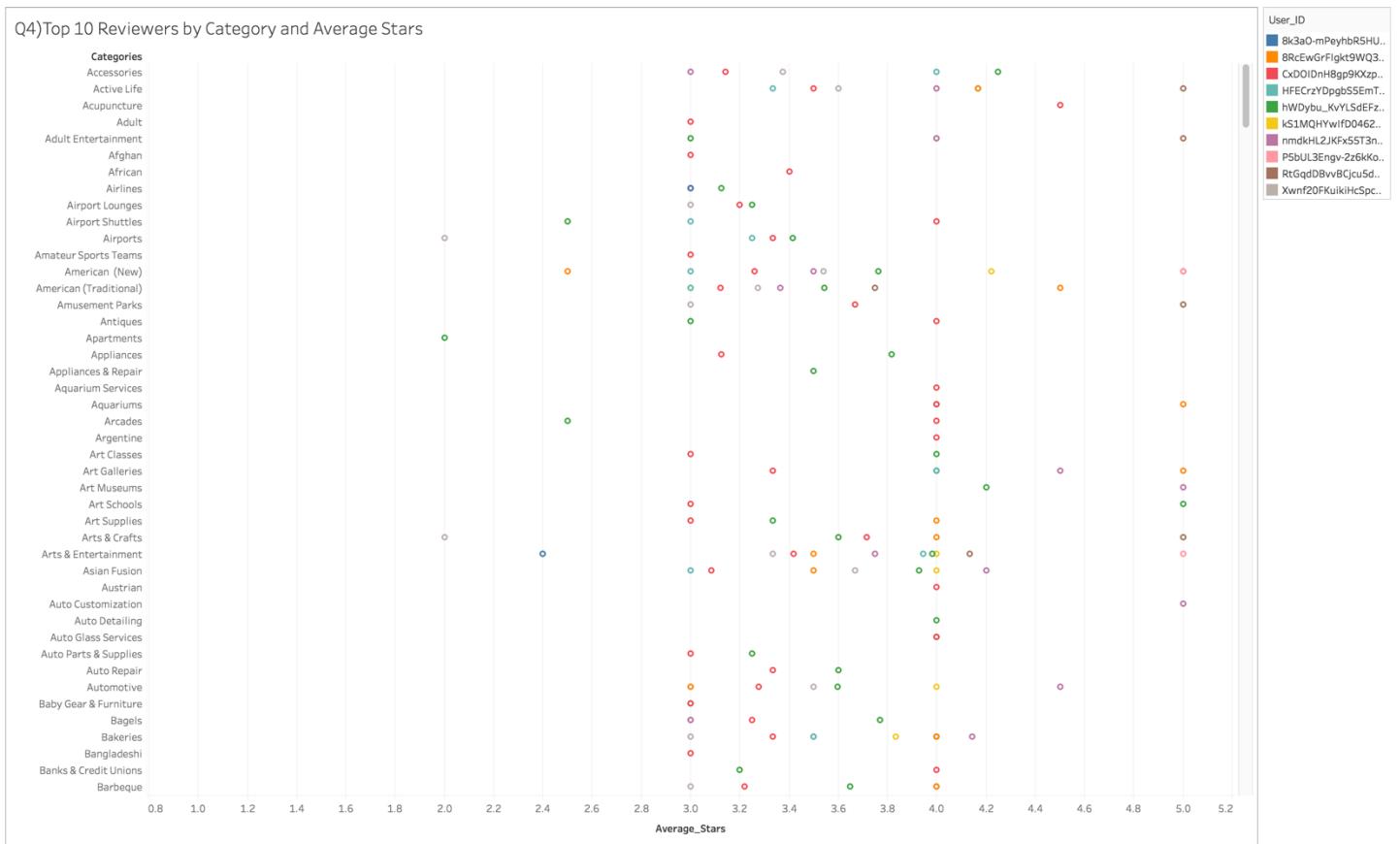
```

Spark Script:

```
1 |from pyspark.sql import SQLContext
2 |sqlContext = SQLContext(sc)
3 |df_business = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_business.json")
4 |df_review = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_review.json")
5 |df_user = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_user.json")
6 |df_review.registerTempTable("reviews")
7 |top10=df_user.sort(df_user.review_count.desc())
8 |top10=top10.limit(10)
9 |top10.registerTempTable("top10")
10 |businesses=df_business.select(pyspark.sql.functions.explode(
11 |    df_business.categories).alias("category"),df_business.business_id,df_business.stars,df_business.city)
12 |businesses.registerTempTable("business")
13 |result=sqlContext.sql("select AVG(reviews.stars) as average_stars, reviews.user_id, category from
14 |    business, reviews, top10 where top10.user_id=reviews.user_id and reviews.business_id=business.business_id
15 |    group by reviews.user_id, category order by reviews.user_id ASC")
16 |result.write.mode('append').json("/user/cloudera/output_spark2/Q4spark.json")
```

Analysis & Visualizations:

In the Visualization, I have shown Average Star rating for Top 10 Users with high review count, I have shown them in colors of scatter plot within different cities.



5. For the top 10 and bottom 10 food business near UWM (in terms of stars), summarize star rating for reviews in January through May.

Pig Script:

```

1 SET elephantbird.jsonloader.nestedLoad 'true';
2
3 A = LOAD './yelp_academic_dataset_business.json' USING com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad=true') AS (yelp: map[]);
4
5 business = FOREACH A GENERATE yelp#`categories` as categories, (int)yelp#`review_count` as review_count, yelp#`city` as city,
6 yelp#`state` as state, (float)yelp#`latitude` as latitude, (float)yelp#`longitude` as longitude, (float)yelp#`stars` as stars;
7
8 coordinates_business = FILTER business BY (latitude<43.221261) AND (latitude>42.931739) AND (longitude<-89.214450) AND
9 (longitude>-89.610310);
10
11 joined = JOIN coordinates_business by business_id, review by business_id;
12
13 businesses = FOREACH joined GENERATE FLATTEN(categories), coordinates_business::stars, coordinates_business::date;
14
15 filter_food = FILTER businesses BY categories MATCHES `.*(Food|food).*`;
16
17 ordered_1 = ORDER filter_food BY stars DESC;
18
19 top_ten = LIMIT ordered_1 10;
20
21 ordered_2 = ORDER filter_food BY stars ASC;
22
23 bottom_ten = LIMIT ordered_2 10;
24
25 top_bottom_20 = UNION top_ten, bottom_ten;
26
27 top_bottom = FOREACH top_bottom_20 GENERATE business.date, stars, categories;
28
29 joined2 = FOREACH top_bottom GENERATE SUBSTRING(date,5,7) as month, stars;
30
31 filter_food2 = FILTER joined2 BY (month matches '01|02|03|04|05|06');
32
33 result_ordered = ORDER filter_food2 BY stars DESC;
34
35 store result_ordered into './FinalQ57';

```

Spark Script:

```

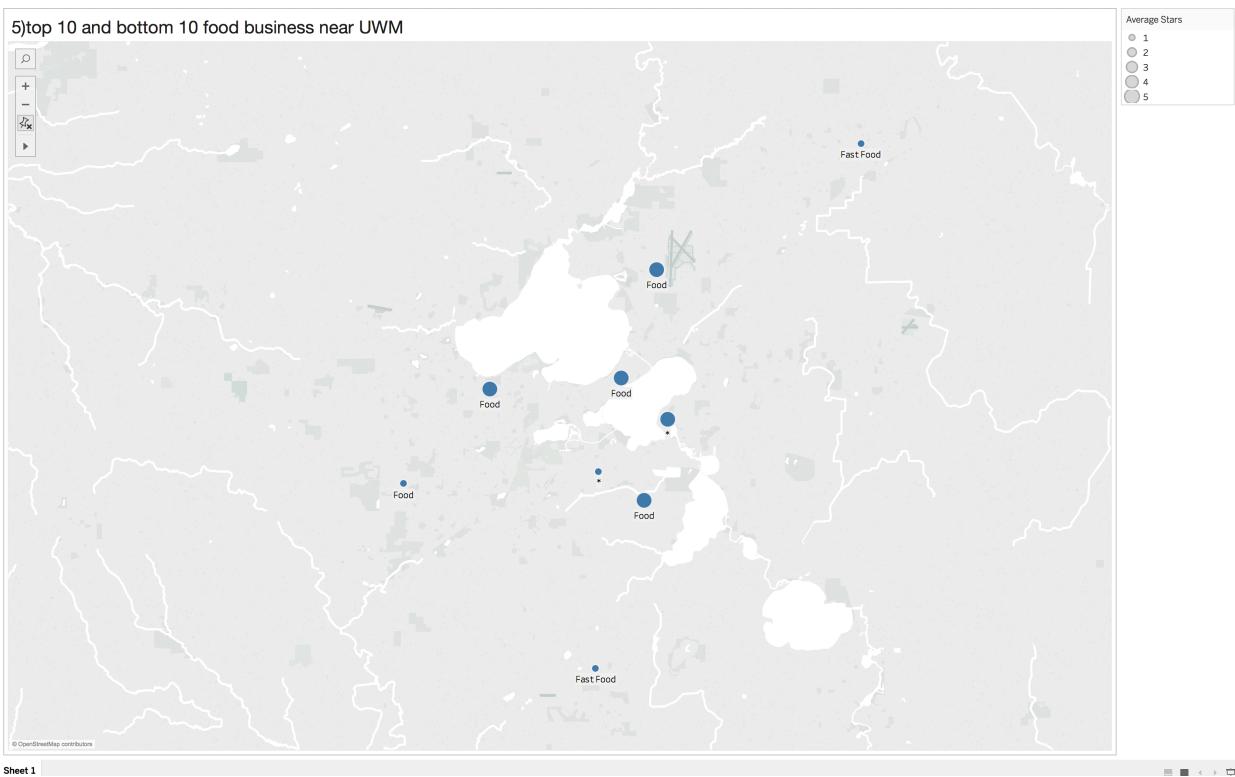
1 from pyspark.sql import SQLContext
2 sqlContext = SQLContext(sc)
3 df_business = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_business.json")
4 df_review = sqlContext.read.json("/user/cloudera/yelp_academic_dataset_review.json")
5 df_business.registerTempTable("business")
6 df_review.registerTempTable("reviews")
7 area=sqlContext.sql("select business_id,stars from business where latitude > 42.931739 and latitude < 43.221261 and longitude > -89.610310 and longitude < -89.214450")
8 top10=area.sort(area.stars.desc())
9 top10=top10.limit(10)
10 top10.registerTempTable("top10")
11 bottom10=area.sort(area.stars.asc())
12 bottom10=bottom10.limit(10)
13 bottom10.registerTempTable("bottom10")
14 reviewmonth=df_review.select(pyspark.sql.functions.year(df_review.date).alias('year'),pyspark.sql.functions.month(df_review.date).alias('month'),df_review.business_id,df_review.stars)
15 reviewmonth.registerTempTable("reviewmonth")
16 resulttop10=sqlContext.sql("select AVG(reviewmonth.stars) as average_stars,reviewmonth.business_id,reviewmonth.year,reviewmonth.month from reviewmonth,top10 where reviewmonth.business_id=top10.business_id group by reviewmonth.business_id,year,month")
17 resultbottom10=sqlContext.sql("select AVG(reviewmonth.stars) as average_stars,reviewmonth.business_id,reviewmonth.year,reviewmonth.month from reviewmonth,bottom10 where reviewmonth.business_id=bottom10.business_id group by reviewmonth.business_id,year,month")
18 reviewmonth.business_id=sqlContext.sql(SELECT FROM Customers
19 WHERE City LIKE '%es%';
20 resulttop10.write.mode('append').json("/user/cloudera/output_spark2/Q5_top10_spark.json")
21 resultbottom10.write.mode('append').json("/user/cloudera/output_spark/Q5_bottom10_spark.json")
22

```

Assumptions:

- I have considered 10 miles from center of UWM with latitudes (42.931739 to 43.221261) and longitudes (-89.610310 to -89.214450)

Analysis & Visualizations:



	A	B	C	D	E	F	G
1	Food Catego	City	State	Average_Star	Longitude	Latitude	ZipCode
2	Fast Food	Sun Prairie	WI	1	-89.2315	43.1866	53590
3	Fast Food	Oregon	WI	1	-89.38384	42.9422	53575
4	Food	Madison	WI	1	-89.50551	43.05725	53719
5	Fast Food	Sun Prairie	WI	1	-89.229965	43.187935	53590
6	Food	Madison	WI	1	-89.32394	43.118195	53704
7	Fast Food	Madison	WI	1	-89.4074	43.0315	53713
8	Food	Madison	WI	1	-89.4074	43.0315	53713
9	Fast Food	Sun Prairie	WI	1	-89.27047	43.16638	53590
10	Fast Food	Madison	WI	1	-89.39466	43.039513	53713
11	Fast Food	Fitchburg	WI	1	-89.40951	43.029854	53713
12	Food	Madison	WI	5	-89.3812	43.074276	53703
13	Food	Monona	WI	5	-89.32558	43.07011	53716
14	Food Trucks	Monona	WI	5	-89.31221	43.047745	53716
15	Food	Monona	WI	5	-89.31221	43.047745	53716
16	Food	Madison	WI	5	-89.42924	43.074207	53705
17	Food	Madison	WI	5	-89.442566	43.075603	53705
18	Food	Fitchburg	WI	5	-89.47182	43.021145	53711
19	Food	Madison	WI	5	-89.34364	43.102917	53704
20	Food	Madison	WI	5	-89.47064	43.078316	53705
21							