

# Assignment-2



**Course: CS 9223**

**Semester: Spring 2017**

**Advisor: Professor Juan Rodriguez**

**Submitted By: Shivaraj Shankaralinga Nesaragi**

**Net ID: ssn314**

**UniversityID: N1558346**

**Date of Submission: 3/12/2017**

## INITIAL SETUP:

These are the initial things that need to be done before you start writing your Pig Script for yelp dataset.

1. Download the yelp dataset from the Yelp website and untar the file to extract 5 Json files.
2. Upload these 5 json files in the hpc. I uploaded these files in the path /user/ssn314/yelp/ folder.
3. Download the following jars: **elephant-bird-core-4.15-RC1.jar**, **elephant-bird-hadoop-compat-4.15-RC1.jar**, **elephant-bird-pig-4.15-RC1.jar**.
4. Upload these jars to your Pig Script Interface in the hpc. I uploaded these files in the path /user/ssn314/.

**Q1) Summarize the number of reviews by US city, by business category.**

### Steps:

- Open up the pig script terminal and click on New Script. Click on the properties tab and add the 3 elephant jars.
- The first line of Pig Script should say **SET elephantbird.jsonloader.nestedLoad 'true'**.
- After this is done, load the yelp\_academic\_dataset\_business.json into the workspace.
- The next step would be to select the fields city, review count and the flattened categories from the json File loaded in the previous step.
- We then need to group the data got from the previous step by city and categories.
- The next step would be to select the cities, categories and sum of the review count from the group data. This is done in Pig with the help of FOREACH and GENERATE commands.
- Once this is done, we use FILTER command to get the data for US cities.
- Then we ORDER by city in order to get the data in a proper sorted format.
- We then store the data got from previous step in the output file. The path of the output file is /user/ssn314/PigOutput/q1Output.

### Script:

```
SET elephantbird.jsonloader.nestedLoad 'true';
```

```
yelp_business_data = LOAD '/user/ssn314/yelp/yelp_academic_dataset_business.json' USING  
com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad') as (json:map[]);
```

```
Yelp_business_data_category = FOREACH yelp_business_data GENERATE (int)json#'review_count'  
AS review_count, json#'city' as city, FLATTEN(json#'categories') AS categories;
```

```

yelp_business_group= GROUP yelp_business_data_category BY (city,categories);

yelp_business_group_data= FOREACH yelp_business_group GENERATE group.city as
city,group.categories as category,SUM(yelp_business_data_category.review_count) AS
reviewCount;

yelp_business_us_cities= FILTER yelp_business_group_data BY city MATCHES
'.*(Charlotte|Cleveland|Madison|Pittsburgh|Phoenix|Urbana|Las Vegas).*';

yelp_business_order= ORDER yelp_business_us_cities BY city;

STORE yelp_business_order INTO '/user/ssn314/PigOutput/q1Output';

```

### **Spark Script:**

```

import spark.implicits._
import org.apache.spark.sql.functions._
import scala.collection.mutable.WrappedArray

val spark = new org.apache.spark.sql.SQLContext(sc)

val business = spark.read.json("/usr/shivraj/yelp_academic_dataset_business.json")

val b = business.withColumn("category", explode(
    when(col("categories").isNotNull, col("categories"))
    .otherwise(array(lit(null).cast("string"))))
))

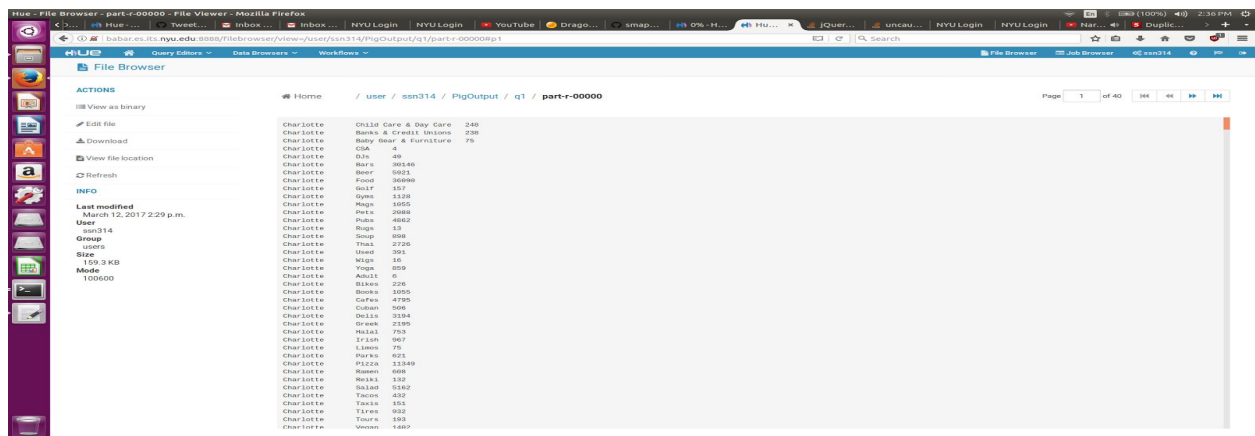
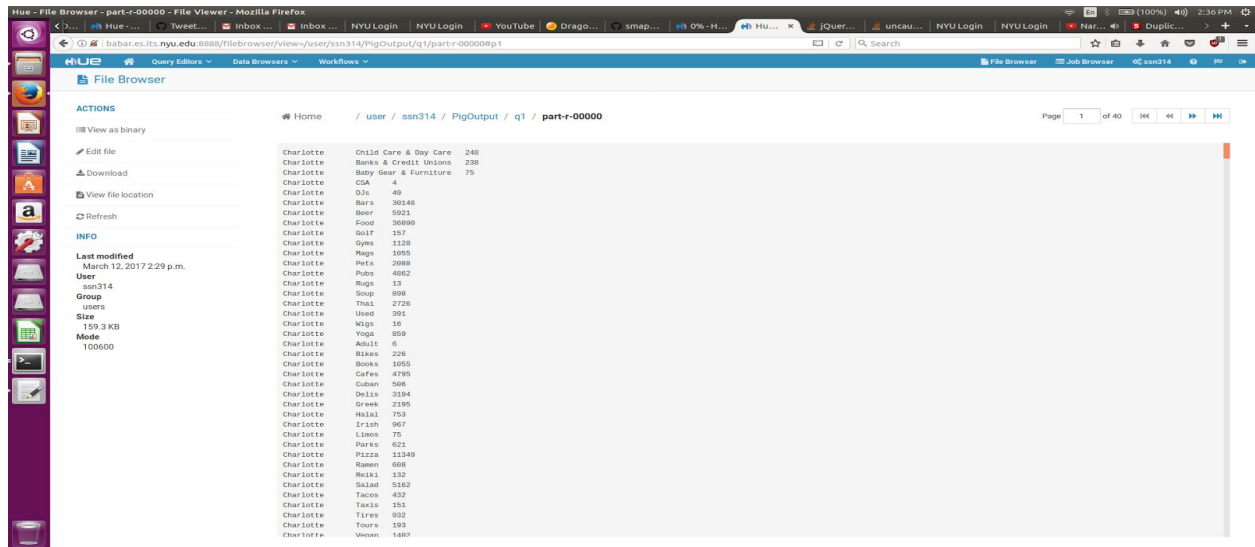
b.registerTempTable("business")

val business_df = spark.sql("select city,category,sum(review_count) from business where
city='Madison' OR city='Pittsburgh' OR city='Charlotte' OR city='Phoenix' OR city='Cleveland' OR
city='Urbana' OR city='Las Vegas' group by city,category Order by city")
business_df.write
    .option("header", "true")
    .csv("/usr/shivraj/q1Spark.csv")

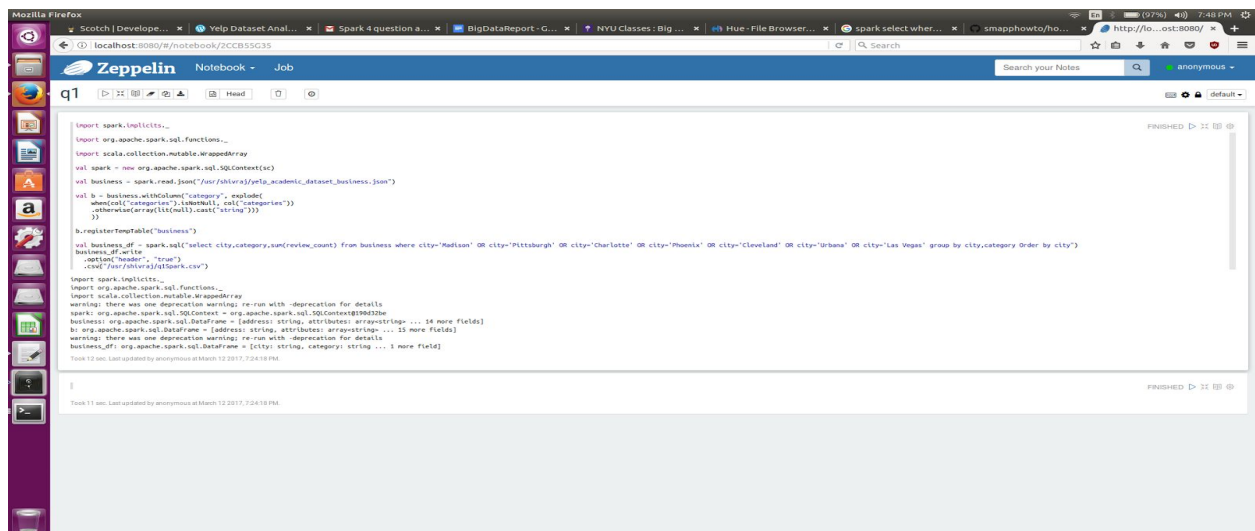
```

### **OUTPUT:**

Pig Script Output:



## Spark Script Output:



**Q2) Rank all *cities* by # of stars descending, for each category.**

**Steps:**

- Open up the pig script terminal and click on New Script. Click on the properties tab and add the 3 elephant jars.
- The first line of Pig Script should say **SET elephantbird.jsonloader.nestedLoad 'true'**.
- After this is done, load the yelp\_academic\_dataset\_business.json into the workspace.
- Select the stars, city and flattened categories from the business Json and store it in a relation.
- Then group the previous relation by categories and city in 'yelp\_business\_group' relation.
- Then generate the fields categories, city and average of stars and store it in **yelp\_business\_group\_data** relation.
- Then sort the relation by category ascending and stars descending and store the relation in **q2FinalOutput.out**.
- The final output contains the cities ranked by stars descending for each category.

**Script:**

```
SET elephantbird.jsonloader.nestedLoad 'true';
```

```
yelp_business_data = LOAD '/user/ssn314/yelp/yelp_academic_dataset_business.json' USING  
com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad') as (json:map[]);
```

```
yelp_business_data_category = FOREACH yelp_business_data GENERATE (float)json#'stars' AS  
stars,json#'city' as city, FLATTEN(json#'categories') AS categories;
```

```
yelp_business_group= GROUP yelp_business_data_category BY (categories,city);
```

```
yelp_business_group_data= FOREACH yelp_business_group GENERATE group.categories as  
category,group.city AS city, AVG(yelp_business_data_category.stars) AS stars;
```

```
yelp_data_order= ORDER yelp_business_group_data BY category ASC,stars DESC;
```

```
STORE yelp_data_order INTO '/user/ssn314/PigOutput/q2FinalOutput.out';
```

**Spark Script:**

```
import spark.implicits._
```

```
import org.apache.spark.sql.functions._
```

```
import scala.collection.mutable.WrappedArray
```

```
val spark = new org.apache.spark.sql.SQLContext(sc)
```

```
val business = spark.read.json("/usr/shivraj/yelp_academic_dataset_business.json")
```

```
val b = business.withColumn("category", explode(
    when(col("categories").isNotNull, col("categories"))
    .otherwise(array(lit(null).cast("string"))))
))
```

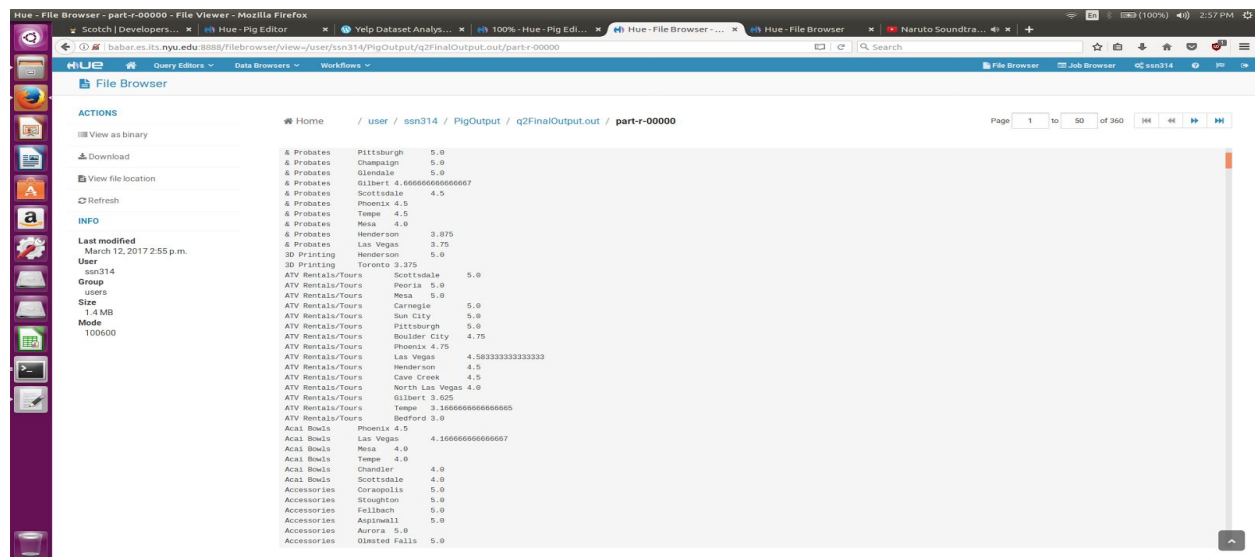
```
b.registerTempTable("business")
```

```
val business_df = spark.sql("select category,city,AVG(stars) as stars from business group by
category,city Order by category,stars DESC")
```

```
business_df.write
    .option("header", "true")
    .csv("/usr/shivraj/q2Spark.csv")
```

## OUTPUT:

### Pig Script Output:



ATV Rentals/Tours	Scottsdale	5.0
ATV Rentals/Tours	Peoria	5.0
ATV Rentals/Tours	Mesa	5.0
ATV Rentals/Tours	Carnegie	5.0
ATV Rentals/Tours	Sun City	5.0
ATV Rentals/Tours	Pittsburgh	5.0
ATV Rentals/Tours	Boulder City	4.75
ATV Rentals/Tours	Phoenix	4.75
ATV Rentals/Tours	Las Vegas	4.583333333333333
ATV Rentals/Tours	Henderson	4.5
ATV Rentals/Tours	Cave Creek	4.5
ATV Rentals/Tours	North Las Vegas	4.0
ATV Rentals/Tours	Gilbert	3.625
ATV Rentals/Tours	Tempe	3.166666666666667
ATV Rentals/Tours	Bedford	3.0
Acad Bowls	Phoenix	4.5
Acad Bowls	Las Vegas	4.166666666666667
Acad Bowls	Mesa	4.0
Acad Bowls	Tempe	4.0
Acad Bowls	Chandler	4.0
Acad Bowls	Scottsdale	4.0
Accessories	Corona	5.0
Accessories	Stoughton	5.0
Accessories	Fellbach	5.0
Accessories	Aspinwall	5.0
Accessories	Aurora	5.0
Accessories	Olathe	5.0

File Browser

ACTIONS

View as binary

Download

View file location

Refresh

INFO

Last modified: March 12, 2017 2:55 p.m.

User: ssn314

Group: users

Size: 1.4 MB

Mode: 100600

Yoga	Brossard	3.0
Yoga	Brooklyn	3.0
Yoga	Lynhurst	3.0
Yoga	East Owellimbury	3.0
Yoga	Wenduridge	3.0
Yoga	Stowe	3.0
Yoga	Brampton	2.9285714285714284
Yoga	Pleasant Hills	2.5
Yoga	Bridgeville	2.5
Yoga	Ajax	2.5
Yoga	Cuyahoga Falls	2.5
Yoga	Harrisburg	2.0
Ziplining	Montreal	4.0
Ziplining	Las Vegas	3.5
Ziplining	Streetsboro	3.0
Zoos	Mesa	5.0
Zoos	Dunfermline	4.5
Zoos	Wenderson	4.5
Zoos	Livingston	4.5
Zoos	Pittsburgh	4.25
Zoos	Toronto	4.1
Zoos	Edinburgh	4.0
Zoos	Scottsdale	4.0
Zoos	Champaign	4.0
Zoos	Phoenix	4.0
Zoos	Madison	4.0
Zoos	Concord	4.0
Zoos	Esslingen	4.0
Zoos	Las Vegas	3.842857142857143
Zoos	Cave Creek	3.5
Zoos	Blue Diamond	3.5
Zoos	Siouxness	3.5
Zoos	Litchfield Park	3.5
Zoos	Stuttgart	3.375
Zoos	Cleveland	3.25
Zoos	Sainte-Anne-de-Bellevue	3.0
Zoos	Glenhale	3.0
Zoos	Saint-Bernard-de-Lacolle	3.0

## Spark Script Output:

```

import org.apache.spark.sql.functions._
import org.apache.spark.sql.functions._
import scala.collection.mutable.WrappedArray

val spark = new org.apache.spark.sql.SQLContext(sc)
val business = spark.read.json("usr/shivraj/yelp_academic_dataset_business.json")

val b = business.withColumn("category", explode(
  whenNotNull("categories").lithemall(), col("categories"))
  .otherwise(array(lit(null)).cast("string")))
)

b.registerTempTable("business")

val business_df = spark.sql("select category,city,avg(stars) as stars from business group by category,city order by category,stars desc")

business_df.write
.option("header", "true")
.csv("usr/shivraj/output.csv")

import org.apache.spark.sql.functions._
import scala.collection.mutable.WrappedArray

warning: there was one deprecation warning; re-run with -deprecation for details
spark: org.apache.spark.sql.SQLContext = org.apache.spark.sql.SQLContext@19638a7
business: org.apache.spark.sql.DataFrame = [address: string, attributes: array[string] ... 14 more fields]
bi: org.apache.spark.sql.DataFrame = [address: string, attributes: array[string] ... 15 more fields]
warning: there was one deprecation warning; re-run with -deprecation for details
business_df: org.apache.spark.sql.DataFrame = [category: string, city: string ... 1 more field]

Task 16 run. Last updated by anonymous at March 12 2017, 7:28:23 PM.
  
```

**Q3) What is the average rank (# stars) for businesses within 10 miles of the University of Wisconsin - Madison, by type of business?**

## Steps:

- Open up the pig script terminal and click on New Script. Click on the properties tab and add the 3 elephant jars.
- The first line of Pig Script should say **SET elephantbird.jsonloader.nestedLoad 'true'**.
- After this is done, load the yelp\_academic\_dataset\_business.json into the workspace.
- Then we select the stars, latitude, longitude, flattened categories and store in the relation **yelp\_business\_data\_category**.
- Then we need to filter the previous relation based on the bounding box near Wisconsin. This is done using FILTER command and then filtered based on latitude and longitude values.

- Then GROUP the relation by categories and store in the relation **yelp\_Wisconsin\_Category**
- Then from this relation select categories from group and average of the stars.
- Then we ORDER by stars descending to get the categories rated according to stars.

### **Script:**

```
SET elephantbird.jsonloader.nestedLoad 'true';
```

```
yelp_business_data = LOAD '/user/ssn314/yelp/yelp_academic_dataset_business.json' USING
                        com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad') as (json:map[]);
```

```
yelp_business_data_category = FOREACH yelp_business_data GENERATE
                                (float)json#'stars' AS stars,
                                (float)json#'latitude' AS latitude,
                                (float)json#'longitude' AS longitude,
                                FLATTEN(json#'categories') AS categories;
```

```
yelp_Wisconsin= FILTER yelp_business_data_category BY (latitude>42.9083)
                  AND (latitude<43.2417)
                  AND (longitude>-89.5839)
                  AND (longitude<-89.2506);
```

```
yelp_Wisconsin_Category= GROUP yelp_Wisconsin BY categories;
```

```
yelp_Wisconsin_BoundingBox= FOREACH yelp_Wisconsin_Category generate group AS categories,
                                AVG(yelp_Wisconsin.stars) AS stars;
```

```
yelp_Wisconsin_BoundingBox_Ordered= ORDER yelp_Wisconsin_BoundingBox BY stars DESC;
```

```
STORE yelp_Wisconsin_BoundingBox_Ordered INTO '/user/ssn314/PigOutput/q3Final.out';
```

### **SPARK SCRIPT:**

```
import spark.implicits._
```

```
import org.apache.spark.sql.functions._
```

```
import scala.collection.mutable.WrappedArray
```



```
val spark = new org.apache.spark.sql.SQLContext(sc)
```

```
val business = spark.read.json("/usr/shivraj/yelp_academic_dataset_business.json")
```

```
val b = business.withColumn("category", explode(
  when(col("categories").isNotNull, col("categories"))
  .otherwise(array(lit(null).cast("string"))))
))
```

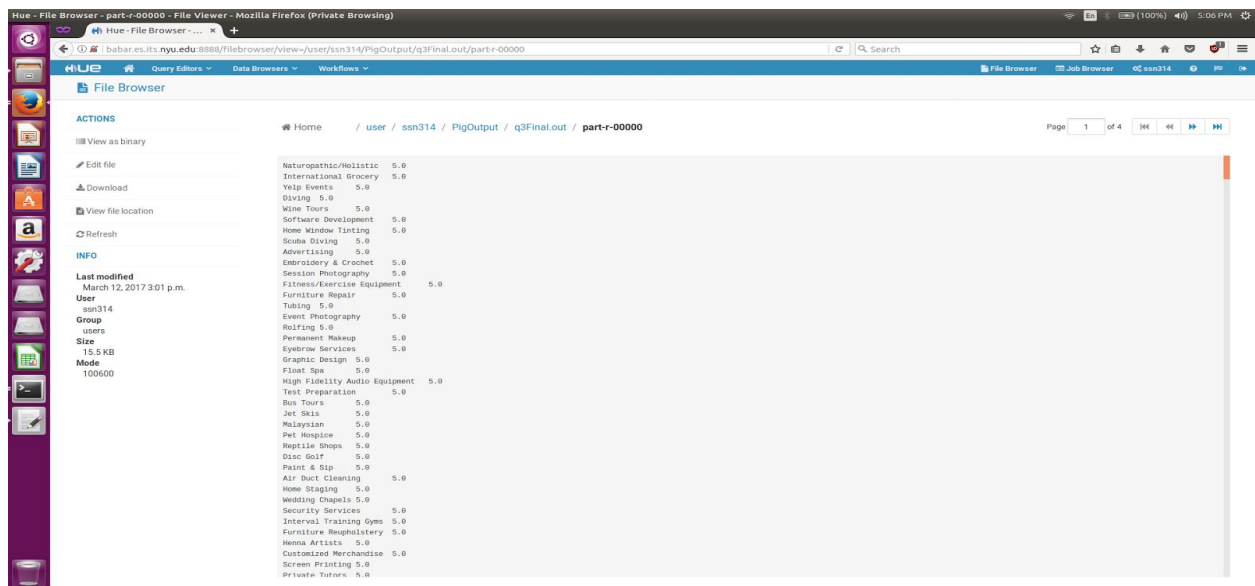
```
b.registerTempTable("business")
```

```
val q3Output=spark.sql("SELECT category,AVG(stars) AS stars from business WHERE
latitude<=43.2467 and latitude>=42.90833 and longitude>=-89.58389 and longitude<=-89.25056
GROUP BY category ORDER BY stars DESC")
```

```
q3Output.write
.option("header", "true")
.csv("/usr/shivraj/q3Spark.csv")
```

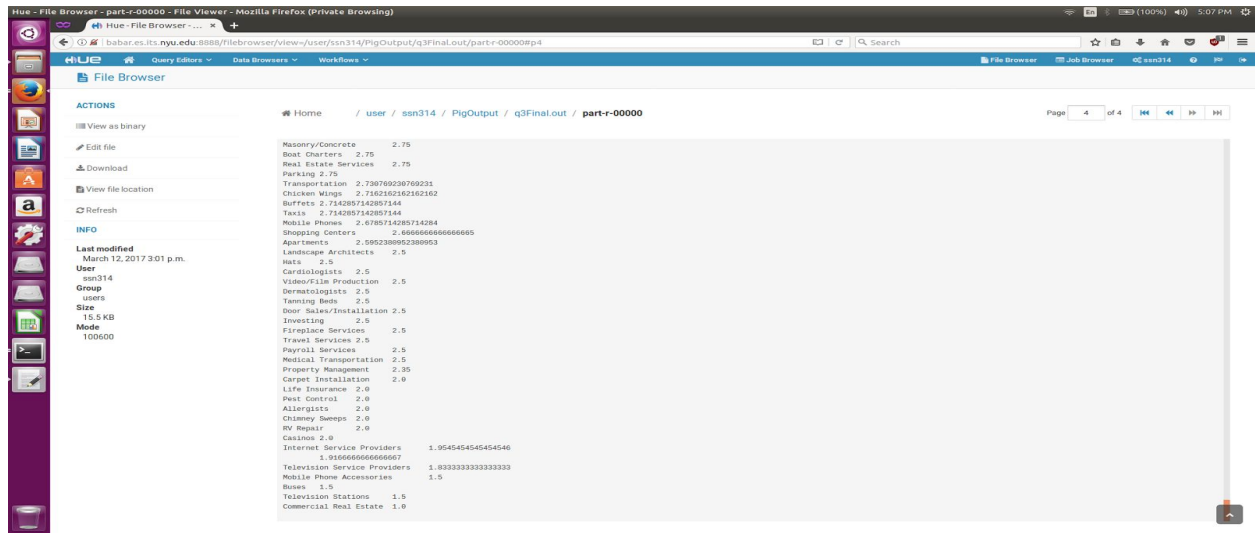
## OUTPUT:

### Pig Script Output:

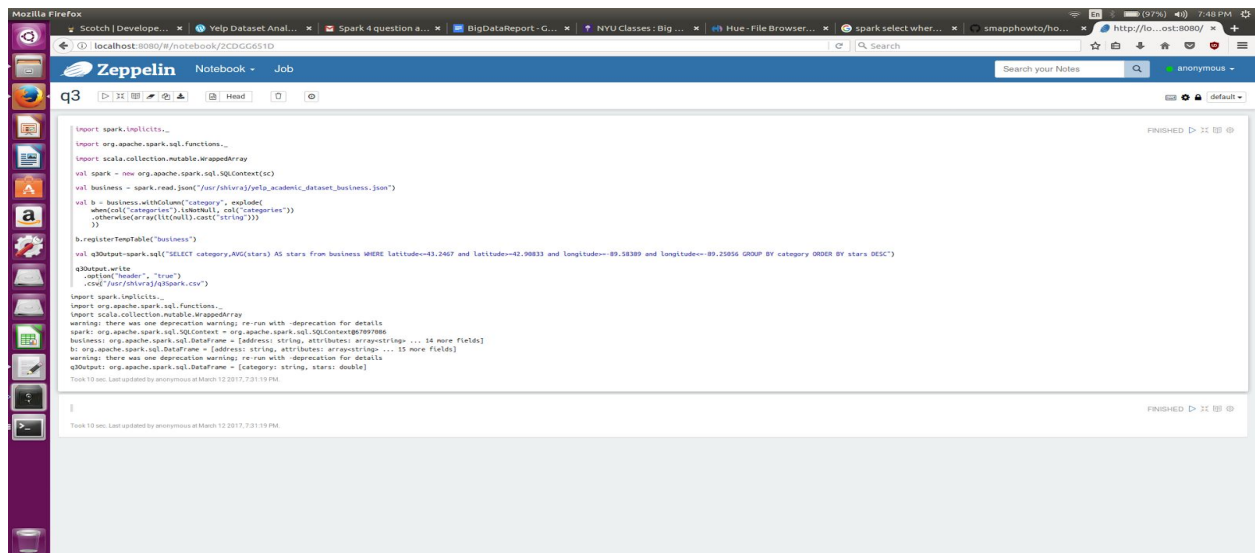


The screenshot shows the Hue File Browser interface. The left sidebar contains a vertical toolbar with icons for various actions like 'View as binary', 'Edit file', 'Download', 'View file location', 'Refresh', and 'Info'. The main panel displays a directory listing for the path '/ user / ssn314 / PigOutput / q3Final.out / part-r-00000'. The listing shows a series of files, each with a name and a rating of 5.0. The files are listed in a single column, with the rating '5.0' appearing to the right of each file name. The files include: Naturopathic/Holistic, International Grocery, Yelp Events, Diving, Wine Tours, Software Development, Home Window Tinting, Scuba Diving, Advertising, Embroidery & Crochet, Session Photography, Fitness/Exercise Equipment, Furniture Repair, Tubing, Event Photography, Rifleing, Permanent Makeup, Eyebrow Services, Graphic Design, Float Spa, High Fidelity Audio Equipment, Test Preparation, Bus Tours, Jet Skis, Malaysian, Pet Hospice, Reptile Shops, Disc Golf, Paint & Sip, Air Duct Cleaning, Home Staging, Wedding Chapels, Security Services, Interval Training Gyms, Furniture Reupholstery, Hanna Artists, Customized Merchandise, Screen Printing, and Private Tutor.

File Name	Rating
Naturopathic/Holistic	5.0
International Grocery	5.0
Yelp Events	5.0
Diving	5.0
Wine Tours	5.0
Software Development	5.0
Home Window Tinting	5.0
Scuba Diving	5.0
Advertising	5.0
Embroidery & Crochet	5.0
Session Photography	5.0
Fitness/Exercise Equipment	5.0
Furniture Repair	5.0
Tubing	5.0
Event Photography	5.0
Rifleing	5.0
Permanent Makeup	5.0
Eyebrow Services	5.0
Graphic Design	5.0
Float Spa	5.0
High Fidelity Audio Equipment	5.0
Test Preparation	5.0
Bus Tours	5.0
Jet Skis	5.0
Malaysian	5.0
Pet Hospice	5.0
Reptile Shops	5.0
Disc Golf	5.0
Paint & Sip	5.0
Air Duct Cleaning	5.0
Home Staging	5.0
Wedding Chapels	5.0
Security Services	5.0
Interval Training Gyms	5.0
Furniture Reupholstery	5.0
Hanna Artists	5.0
Customized Merchandise	5.0
Screen Printing	5.0
Private Tutor	5.0



## Spark Script Output:



**Q4) Rank reviewers by number of reviews. For the top 10 reviewers, show their average number of stars, by category.**

### Steps:

- Open up the pig script terminal and click on New Script. Click on the properties tab and add the 3 elephant jars.
- The first line of Pig Script should say **SET elephantbird.jsonloader.nestedLoad 'true'**.
- After this is done, load the yelp\_academic\_dataset\_user.json into the workspace.
- Select the userId and review count from the userJson and store in the **user\_data** relation.
- Select the top 10 reviewers from then.
- Then load the review Json into the workspace.
- Select the userId ,stars and businessId from the review Json.

- We then join user relation for top 10 reviewers and review Json by userId.
- We then select the userId, businessId and stars from this joined relation and store it in **review\_user\_join** relation..
- The next step is to load the business Json.
- Then we select business Id and flattened categories and store in the relation **categoryID**
- We then join the **categoryID** relation and the **review\_user\_join** relation by businessId.
- Then we select userId,category and stars from this relation and group it by userid and category and store it in **reviewGrouped** relation.
- The last step is to select the userId,category and average of stars for each category and store the relation in the output directory **/user/ssn314/PigOutput/Q4FinalOutput.out**.

### Script:

```
SET elephantbird.jsonloader.nestedLoad 'true';
```

```
userJson = LOAD '/user/ssn314/yelp/yelp_academic_dataset_user.json' USING
            com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad') as (users:map[]);
```

```
user_data = FOREACH userJson GENERATE (chararray)users#'user_id' as user_id,
                                       (int)users#'review_count' as review_count;
```

```
user_ordered_data = ORDER user_data by review_count DESC;
```

```
top_10_reviewers = LIMIT user_ordered_data 10;
```

```
reviewJson = LOAD '/user/ssn314/yelp/yelp_academic_dataset_review.json' USING
              com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad') as
              (review:map[]);
```

```
review_data = FOREACH reviewJson GENERATE (chararray)review#'user_id' as user_id,
                                           (float)review#'stars' as stars,(chararray)review#'business_id' as business_id;
```

```
review_user_join = JOIN review_data BY user_id, top_10_reviewers by user_id;
```

```
review_counter = FOREACH review_user_join GENERATE review_data::user_id as user_id,
            review_data::business_id as business_id,review_data::stars as stars;
```

```
businessJson = LOAD '/user/ssn314/yelp/yelp_academic_dataset_business.json' USING
                com.twitter.elephantbird.pig.load.JsonLoader('-nestedLoad') as (business:map[]);
```

```
categoryID = FOREACH businessJson GENERATE (chararray)business#'business_id' AS business_id,
            FLATTEN(business#'categories') AS category;
```

```

reviewJoined = JOIN categoryID BY business_id, review_counter BY business_id;

reviewFinal = FOREACH reviewJoined GENERATE review_counter::user_id as user_id,
                                         categoryID::category as category, review_counter::stars as stars;

reviewGrouped = GROUP reviewFinal by (user_id,category);

finalOutput = FOREACH reviewGrouped GENERATE group.user_id as user_id,
                                         group.category as category, AVG(reviewFinal.stars) as stars;

STORE top_10_reviewers into '/user/ssn314/PigOutput/Q4_top_ten_reviewers.out';

STORE finalOutput into '/user/ssn314/PigOutput/Q4FinalOutput.out';

```

### **Spark Script:**

```

import spark.implicits._

import org.apache.spark.sql.functions._

import scala.collection.mutable.WrappedArray

val spark = new org.apache.spark.sql.SQLContext(sc)

val business = spark.read.json("/usr/shivraj/yelp_academic_dataset_business.json")

val b = business.withColumn("category", explode(
    when(col("categories").isNotNull, col("categories"))
    .otherwise(array(lit(null).cast("string"))))
))

val user = spark.read.json("/usr/shivraj/yelp_academic_dataset_user.json")

val review = spark.read.json("/usr/shivraj/yelp_academic_dataset_review.json")

b.registerTempTable("business")
user.registerTempTable("user")
review.registerTempTable("review")

val top_10_reviewers=spark.sql("SELECT user_id,review_count FROM user ORDER BY review_count
DESC LIMIT 10")

top_10_reviewers.registerTempTable("top_10_reviewers")

```

```
val userReviewJson=spark.sql("SELECT top.user_id,r.business_id,r.stars FROM top_10_reviewers AS top,review AS r WHERE r.user_id=top.user_id")
```

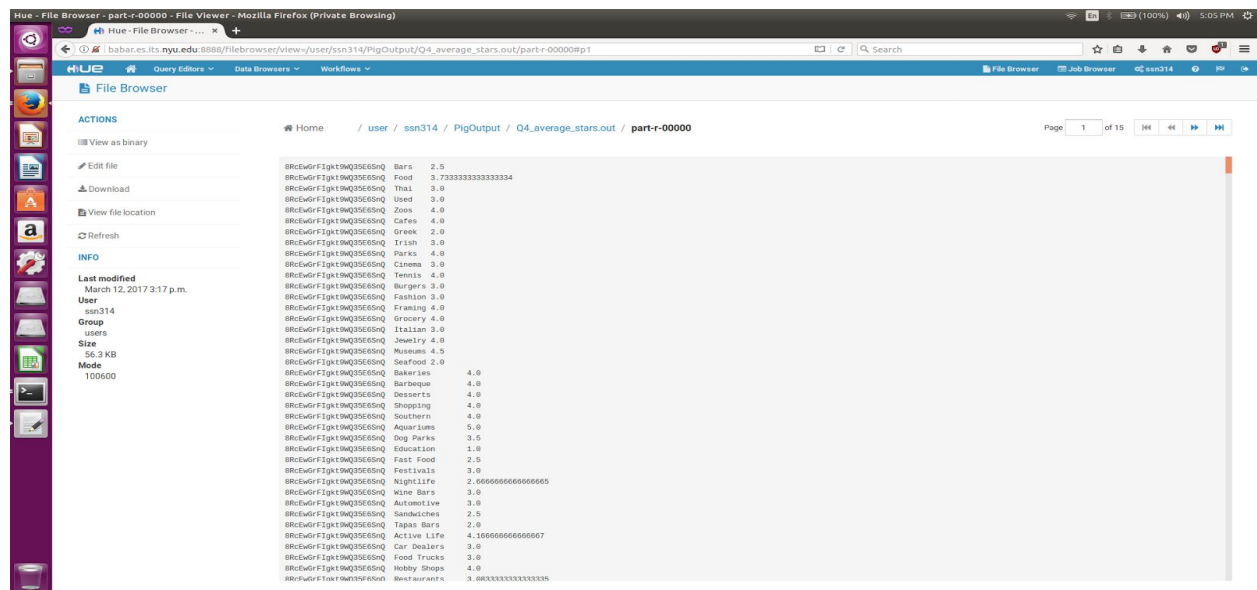
```
userReviewJson.registerTempTable("userReviewJson")
```

```
val q4Final=spark.sql("SELECT u.user_id,avg(u.stars),b.category FROM userReviewJson AS u,business AS b where b.business_id=u.business_id GROUP BY u.user_id,b.category ORDER BY user_id,category")
```

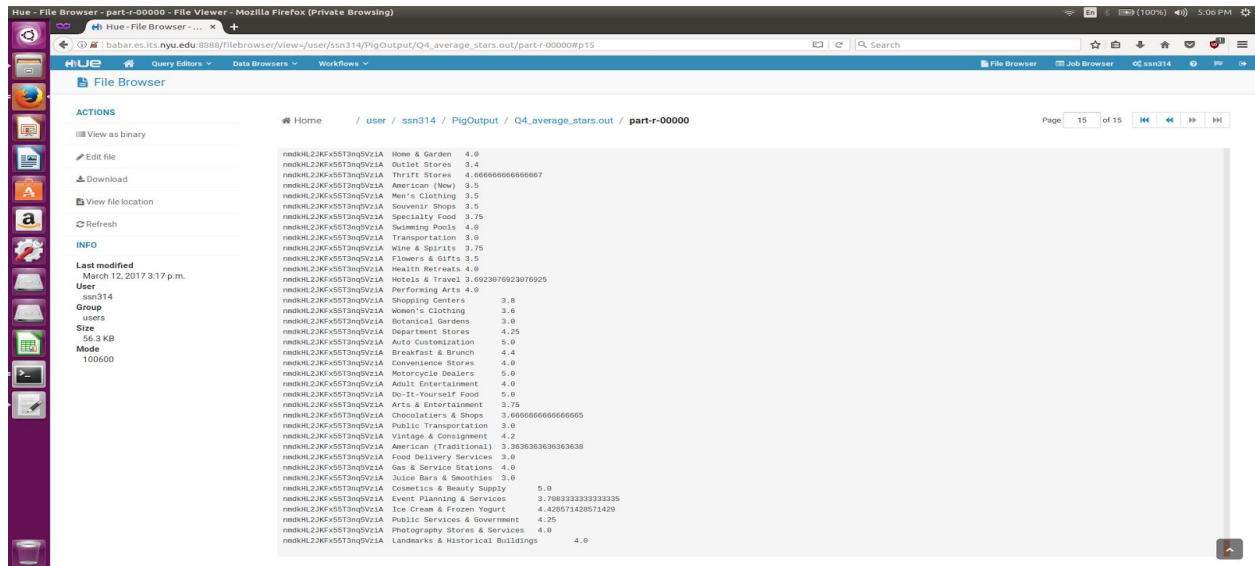
```
q4Final.write
.option("header", "true")
.csv("/usr/shivraj/q4Spark.csv")
```

## OUTPUT:

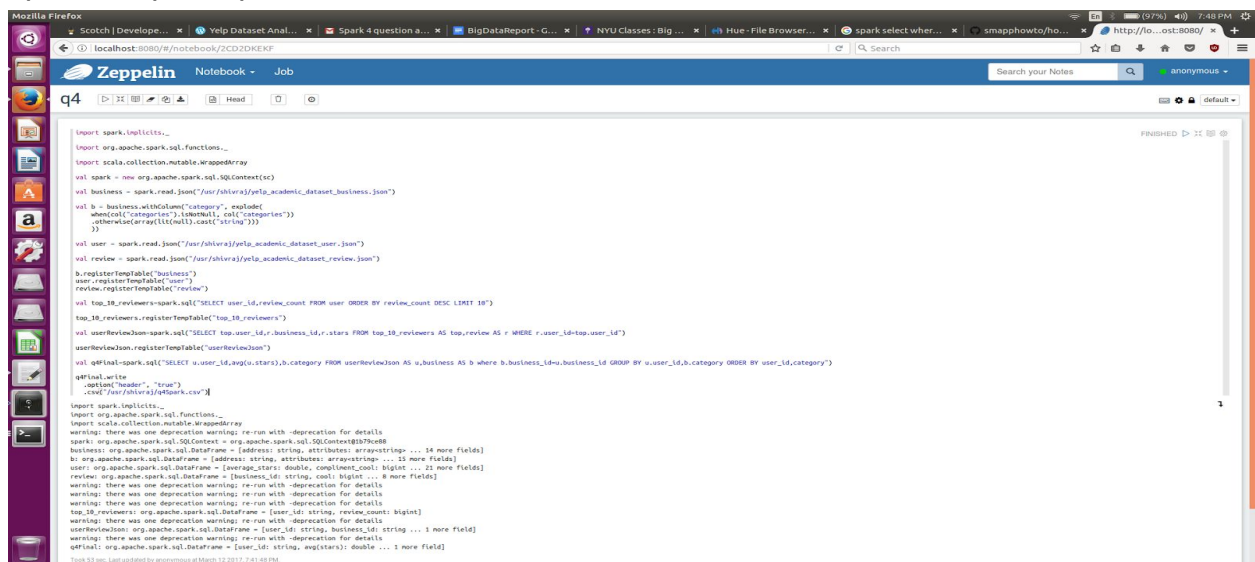
### Pig Script Output:



File Name	Content
part-r-00000	<pre> 8RCEwDrF2gk1t9Mq3SE65nQ Bars 2.5 8RCEwDrF2gk1t9Mq3SE65nQ Food 3.7333333333333334 8RCEwDrF2gk1t9Mq3SE65nQ Thai 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Used 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Zoes 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Cafes 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Greek 2.0 8RCEwDrF2gk1t9Mq3SE65nQ Irish 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Parks 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Cinema 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Tennis 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Burgers 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Fashion 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Framing 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Grocery 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Italian 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Jewelry 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Museums 4.5 8RCEwDrF2gk1t9Mq3SE65nQ Seafood 2.0 8RCEwDrF2gk1t9Mq3SE65nQ Bakeries 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Baroque 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Desserts 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Shopping 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Southern 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Aquariums 5.0 8RCEwDrF2gk1t9Mq3SE65nQ Dog Parks 3.5 8RCEwDrF2gk1t9Mq3SE65nQ Education 1.0 8RCEwDrF2gk1t9Mq3SE65nQ Fast Food 2.5 8RCEwDrF2gk1t9Mq3SE65nQ Festivals 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Night Life 2.6666666666666665 8RCEwDrF2gk1t9Mq3SE65nQ Wine Bars 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Automotive 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Sandwiches 2.5 8RCEwDrF2gk1t9Mq3SE65nQ Tapas Bars 2.0 8RCEwDrF2gk1t9Mq3SE65nQ Active Life 4.166666666666667 8RCEwDrF2gk1t9Mq3SE65nQ Car Dealers 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Food Trucks 3.0 8RCEwDrF2gk1t9Mq3SE65nQ Hobby Shops 4.0 8RCEwDrF2gk1t9Mq3SE65nQ Restaurants 3.0000000000000000 </pre>



## Spark Script Output:



**Q5) For the top 10 and bottom 10 food business near UWM (in terms of stars), summarize star rating for reviews in January through May.**

## Steps:

- Open up the pig script terminal and click on New Script. Click on the properties tab and add the 3 elephant jars.
- The first line of Pig Script should say **SET elephantbird.jsonloader.nestedLoad 'true'**.
- After this is done, load the yelp\_academic\_dataset\_business.json and yelp\_academic\_dataset\_review.json into the workspace.
- We select the businessId,latitude,stars,longitude and flattened categories and store it in the **yelp\_business\_data\_category**.

- We select the businessId, reviewId,userId,stars and date and store it in **yelp\_review\_category** relation.
- The next step is to select the reviews which are in the month range of January to May.
- We then select the Food category restaurants that are near Wisconsin.
- Then from that we order the restaurants by stars and take the top 10 and bottom 10 restaurants from this relation.
- We then make the UNION of the top 10 and bottom 10 relations and store it in a relation **yelp\_top\_bottom\_10**.
- Then we join the **yelp\_top\_bottom\_10** and **yelp\_review\_date** by businessId.
- We then select the businessId,reviewId,userId,stars,date.
- We then group by businessID and store it in the relation **yelp\_biz\_group**.
- From the **yelp\_biz\_group** we select businessId, average of stars and date and store the output in **q5Final.out**

### **Script:**

```
SET elephantbird.jsonloader.nestedLoad 'true';
```

```
yelp_business_data = LOAD '/user/ssn314/yelp/yelp_academic_dataset_business.json' USING
com.twitter.elephantbird.pig.load.JsonLoader
                        ('-nestedLoad') as (business:map[]);
```

```
yelp_review_data = LOAD '/user/ssn314/yelp/yelp_academic_dataset_review.json' USING
com.twitter.elephantbird.pig.load.JsonLoader
                        ('-nestedLoad') as (review:map[]);
```

```
yelp_business_data_category = FOREACH yelp_business_data GENERATE
                                (chararray)business#'business_id' AS businessID,
                                (float)business#'latitude' AS latitude,
                                (float)business#'stars' AS stars,
                                (float)business#'longitude' AS longitude,
                                FLATTEN(business#'categories') AS categories;
```

```
yelp_review_category = FOREACH yelp_review_data GENERATE
                        (chararray)review#'business_id' AS businessID,
                        (chararray)review#'review_id' AS reviewID,
                        (chararray)review#'user_id' AS userID,
                        (float)review#'stars' AS stars,
                        (chararray)review#'date' AS date;
```

```
yelp_review_date= FILTER yelp_review_category BY (SUBSTRING(date,5,7) MATCHES
'.*(01|02|03|04|05).*');
```

```

yelp_Wisconsin= FILTER yelp_business_data_category BY (latitude>42.9083)
                AND (latitude<43.2417)
                AND (longitude>-89.5839)
                AND (longitude<-89.2506);

yelp_Wisconsin_Food_Filter= FILTER yelp_Wisconsin BY categories MATCHES '.*Food*.';

yelp_Wisconsin_BoundingBox_desc= ORDER yelp_Wisconsin_Food_Filter BY stars DESC;

yelp_top_10_Wisconsin= LIMIT yelp_Wisconsin_BoundingBox_desc 10;

yelp_Wisconsin_BoundingBox_asc= ORDER yelp_Wisconsin_Food_Filter BY stars ASC;

yelp_bottom_10_Wisconsin= LIMIT yelp_Wisconsin_BoundingBox_asc 10;

yelp_top_bottom_10= UNION yelp_top_10_Wisconsin,yelp_bottom_10_Wisconsin;

STORE yelp_top_bottom_10 INTO '/user/ssn314/PigOutput/q5TopBottomUnion';

yelp_joined_data= JOIN yelp_top_bottom_10 BY businessID,yelp_review_date BY businessID;

STORE yelp_joined_data INTO '/user/ssn314/PigOutput/q5BizIdJoined';

yelp_joined_review_data= FOREACH yelp_joined_data GENERATE yelp_top_bottom_10::businessID
AS bizID,reviewID AS reviewID, userID AS userID,yelp_review_date::stars AS stars,date AS date;

yelp_biz_group= GROUP yelp_joined_review_data BY bizID;

STORE yelp_biz_group INTO '/user/ssn314/PigOutput/q5BizGrouped';

yelp_final_data= FOREACH yelp_biz_group GENERATE group AS bizID,
AVG(yelp_joined_review_data.stars) AS stars,yelp_joined_review_data.date AS date;

STORE yelp_final_data INTO '/user/ssn314/PigOutput/q5Final.out';

```

## **OUTPUT:**

Pig Script Output:



