

# Analysis of Yelp Business Intelligence Data

We will analyze a subset of Yelp's business, reviews and user data. This dataset comes to us from [Kaggle](#) although we have taken steps to pull this data into a public s3 bucket: `s3://sta9760-yelpdataset/yelp-light/*business.json`

## Installation and Initial Setup

Begin by installing the necessary libraries that you may need to conduct your analysis. At the very least, you must install `pandas` and `matplotlib`

In [2]: `%%info`

```
Current session configs: {'conf': {'spark.pyspark.python': 'python3', 'spark.pyspark.virtualenv.enabled': 'true',
'spark.pyspark.virtualenv.type': 'native', 'spark.pyspark.virtualenv.bin.path': '/usr/bin/virtualenv'}, 'kind':
'pyspark'}
No active sessions.
```

In [3]: `sc.install_pypi_package("matplotlib==3.2.1")`  
`sc.install_pypi_package("pandas==1.0.3")`  
`sc.install_pypi_package("scipy==1.7.0")`  
`sc.install_pypi_package("seaborn==0.10.0")`

Starting Spark application

ID	YARN Application ID	Kind	State	Spark UI	Driver log	Current session?
2	application_1638424444107_0003	pyspark	idle			✓

SparkSession available as 'spark'.

Collecting matplotlib==3.2.1

Using cached [https://files.pythonhosted.org/packages/b2/c2/71fcf957710f3ba1f09088b35776a799ba7dd95f7c2b195ec800933b276b/matplotlib-3.2.1-cp37-cp37m-manylinux1\\_x86\\_64.whl](https://files.pythonhosted.org/packages/b2/c2/71fcf957710f3ba1f09088b35776a799ba7dd95f7c2b195ec800933b276b/matplotlib-3.2.1-cp37-cp37m-manylinux1_x86_64.whl)

Collecting python-dateutil>=2.1 (from matplotlib==3.2.1)

Using cached [https://files.pythonhosted.org/packages/36/7a/87837f39d0296e723bb9b62bbb257d0355c7f6128853c78955f57342a56d/python\\_dateutil-2.8.2-py2.py3-none-any.whl](https://files.pythonhosted.org/packages/36/7a/87837f39d0296e723bb9b62bbb257d0355c7f6128853c78955f57342a56d/python_dateutil-2.8.2-py2.py3-none-any.whl)

Collecting pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.1 (from matplotlib==3.2.1)

Using cached <https://files.pythonhosted.org/packages/a0/34/895006117f6fce0b4de045c87e154ee4a20c68ec0a4c9a36d900888fb6bc/pyparsing-3.0.6-py3-none-any.whl>

Collecting cycler>=0.10 (from matplotlib==3.2.1)

Using cached <https://files.pythonhosted.org/packages/5c/f9/695d6bedebd747e5eb0fe8fad57b72fdf25411273a39791cde838d5a8f51/cyclar-0.11.0-py3-none-any.whl>

Requirement already satisfied: numpy>=1.11 in /usr/local/lib64/python3.7/site-packages (from matplotlib==3.2.1)

Collecting kiwisolver>=1.0.1 (from matplotlib==3.2.1)

Using cached [https://files.pythonhosted.org/packages/09/6b/6e567cb2e86d4e5939a9233f8734e26021b6a9c1bc4b1edccba236a84cc2/kiwisolver-1.3.2-cp37-cp37m-manylinux\\_2\\_5\\_x86\\_64.manylinux1\\_x86\\_64.whl](https://files.pythonhosted.org/packages/09/6b/6e567cb2e86d4e5939a9233f8734e26021b6a9c1bc4b1edccba236a84cc2/kiwisolver-1.3.2-cp37-cp37m-manylinux_2_5_x86_64.manylinux1_x86_64.whl)

Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/site-packages (from python-dateutil>=2.1->matplotlib==3.2.1)

```
Installing collected packages: python-dateutil, pyparsing, cycler, kiwisolver, matplotlib
Successfully installed cycler-0.11.0 kiwisolver-1.3.2 matplotlib-3.2.1 pyparsing-3.0.6 python-dateutil-2.8.2
```

```
Collecting pandas==1.0.3
```

```
Using cached https://files.pythonhosted.org/packages/4a/6a/94b219b8ea0f2d580169e85ed1edc0163743f55aaeca8a44c2e8fc1e344e/pandas-1.0.3-cp37-cp37m-manylinux1_x86_64.whl
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.7/site-packages (from pandas==1.0.3)
Requirement already satisfied: numpy>=1.13.3 in /usr/local/lib64/python3.7/site-packages (from pandas==1.0.3)
Requirement already satisfied: python-dateutil>=2.6.1 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from pandas==1.0.3)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/site-packages (from python-dateutil>=2.6.1->pandas==1.0.3)
Installing collected packages: pandas
Successfully installed pandas-1.0.3
```

```
Collecting scipy==1.7.0
```

```
Using cached https://files.pythonhosted.org/packages/b2/85/b00f13b52d079b5625e1a12330fc6453c947a482ff667a907c7bc60ed220/scipy-1.7.0-cp37-cp37m-manylinux_2_5_x86_64.manylinux1_x86_64.whl
Requirement already satisfied: numpy<1.23.0,>=1.16.5 in /usr/local/lib64/python3.7/site-packages (from scipy==1.7.0)
Installing collected packages: scipy
Successfully installed scipy-1.7.0
```

```
Collecting seaborn==0.10.0
```

```
Using cached https://files.pythonhosted.org/packages/70/bd/5e6bf595fe6ee0f257ae49336dd180768c1ed3d7c7155b2fdf894c1c808a/seaborn-0.10.0-py3-none-any.whl
Requirement already satisfied: pandas>=0.22.0 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from seaborn==0.10.0)
Requirement already satisfied: numpy>=1.13.3 in /usr/local/lib64/python3.7/site-packages (from seaborn==0.10.0)
Requirement already satisfied: scipy>=1.0.1 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from seaborn==0.10.0)
Requirement already satisfied: matplotlib>=2.1.2 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from seaborn==0.10.0)
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.7/site-packages (from pandas>=0.22.0->seaborn==0.10.0)
Requirement already satisfied: python-dateutil>=2.6.1 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from pandas>=0.22.0->seaborn==0.10.0)
Requirement already satisfied: pyparsing!=2.0.4,!>=2.1.2,!>=2.1.6,>=2.0.1 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from matplotlib>=2.1.2->seaborn==0.10.0)
Requirement already satisfied: cycler>=0.10 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from matplotlib>=2.1.2->seaborn==0.10.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /mnt/tmp/1638480988030-0/lib/python3.7/site-packages (from matplotlib>=2.1.2->seaborn==0.10.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/site-packages (from python-dateutil>=2.6.1->pandas>=0.22.0->seaborn==0.10.0)
Installing collected packages: seaborn
Successfully installed seaborn-0.10.0
```

## Importing

Now, import the installed packages from the previous block below.

```
In [4]: import pandas as pd
import matplotlib.pyplot as plt
import scipy
import seaborn as sns
```

# Loading Data

We are finally ready to load data. Using `spark` load the data from S3 into a `dataframe` object that we can manipulate further down in our analysis.

```
In [5]: business_df = spark.read.json('s3://jm-p2-bucket/yelp_academic_dataset_business.json')
```

## Overview of Data

Display the number of rows and columns in our dataset.

```
In [6]: print(f'Columns: {len(business_df.dtypes)} | Rows: {business_df.count()}')
```

Columns: 14 | Rows: 160585

Display the DataFrame schema below.

```
In [7]: business_df.printSchema()
```

```
root
|-- address: string (nullable = true)
|-- attributes: struct (nullable = true)
|   |-- AcceptsInsurance: string (nullable = true)
|   |-- AgesAllowed: string (nullable = true)
|   |-- Alcohol: string (nullable = true)
|   |-- Ambience: string (nullable = true)
|   |-- BYOB: string (nullable = true)
|   |-- BYOBCorkage: string (nullable = true)
|   |-- BestNights: string (nullable = true)
|   |-- BikeParking: string (nullable = true)
|   |-- BusinessAcceptsBitcoin: string (nullable = true)
|   |-- BusinessAcceptsCreditCards: string (nullable = true)
|   |-- BusinessParking: string (nullable = true)
|   |-- ByAppointmentOnly: string (nullable = true)
|   |-- Caters: string (nullable = true)
|   |-- CoatCheck: string (nullable = true)
|   |-- Corkage: string (nullable = true)
|   |-- DietaryRestrictions: string (nullable = true)
|   |-- DogsAllowed: string (nullable = true)
|   |-- DriveThru: string (nullable = true)
|   |-- GoodForDancing: string (nullable = true)
|   |-- GoodForKids: string (nullable = true)
|   |-- GoodForMeal: string (nullable = true)
|   |-- HairSpecializesIn: string (nullable = true)
|   |-- HappyHour: string (nullable = true)
|   |-- HasTV: string (nullable = true)
|   |-- Music: string (nullable = true)
```

```

|      |-- NoiseLevel: string (nullable = true)
|      |-- Open24Hours: string (nullable = true)
|      |-- OutdoorSeating: string (nullable = true)
|      |-- RestaurantsAttire: string (nullable = true)
|      |-- RestaurantsCounterService: string (nullable = true)
|      |-- RestaurantsDelivery: string (nullable = true)
|      |-- RestaurantsGoodForGroups: string (nullable = true)
|      |-- RestaurantsPriceRange2: string (nullable = true)
|      |-- RestaurantsReservations: string (nullable = true)
|      |-- RestaurantsTableService: string (nullable = true)
|      |-- RestaurantsTakeOut: string (nullable = true)
|      |-- Smoking: string (nullable = true)
|      |-- WheelchairAccessible: string (nullable = true)
|      |-- WiFi: string (nullable = true)
|-- business_id: string (nullable = true)
|-- categories: string (nullable = true)
|-- city: string (nullable = true)
|-- hours: struct (nullable = true)
|   |-- Friday: string (nullable = true)
|   |-- Monday: string (nullable = true)
|   |-- Saturday: string (nullable = true)
|   |-- Sunday: string (nullable = true)
|   |-- Thursday: string (nullable = true)
|   |-- Tuesday: string (nullable = true)
|   |-- Wednesday: string (nullable = true)
|-- is_open: long (nullable = true)
|-- latitude: double (nullable = true)
|-- longitude: double (nullable = true)
|-- name: string (nullable = true)
|-- postal_code: string (nullable = true)
|-- review_count: long (nullable = true)
|-- stars: double (nullable = true)
|-- state: string (nullable = true)

```

Display the first 5 rows with the following columns:

- business\_id
- name
- city
- state
- categories

```
In [8]: business_df.select('business_id', 'name', 'city', 'state', 'categories').show(5)
```

business_id	name	city	state	categories
6iYb2HFDywm3zjuRg...	Oskar Blues Taproom	Boulder	CO	Gastropubs, Food,...
tCbdrRPZA0oiIYSmH...	Flying Elephants ...	Portland	OR	Salad, Soup, Sand...
bvN78f1M8NLprQ1a1...	The Reclaimory	Portland	OR	Antiques, Fashion...

```
|oaepsyvc0J17qwi8c...|      Great Clips|Orange City|      FL|Beauty & Spas, Ha...|
|PE9uqAjdW0E4-8mjG...|      Crossfit Terminus|      Atlanta|      GA|Gyms, Active Life...|
+-----+-----+-----+-----+-----+
only showing top 5 rows
```

# Analyzing Categories

Let's now answer this question: **how many unique categories are represented in this dataset?**

Essentially, we have the categories per business as a list - this is useful to quickly see what each business might be represented as but it is difficult to easily answer questions such as:

- How many businesses are categorized as `Active Life` , for instance
- What are the top 20 most popular categories available?

## Association Table

We need to "break out" these categories from the business ids? One common approach to take is to build an association table mapping a single business id multiple times to each distinct category.

For instance, given the following:

business_id	categories
abcd123	a,b,c

We would like to derive something like:

business_id	category
abcd123	a
abcd123	b
abcd123	c

What this does is allow us to then perform a myriad of rollups and other analysis on this association table which can aid us in answering the questions asked above.

Implement the code necessary to derive the table described from your original yelp dataframe.

```
In [9]: from pyspark.sql.functions import split, explode

        business_assoc_table = business_df.withColumn('categories', explode(split('categories',', ')))
```

Display the first 5 rows of your association table below.

```
In [10]: business_assoc_table.select('business_id', 'categories').show(5)
```

```
+-----+-----+
|      business_id| categories|
+-----+-----+
|6iYb2HFDywm3zjuRg...| Gastropubs|
|6iYb2HFDywm3zjuRg...|      Food|
|6iYb2HFDywm3zjuRg...|Beer Gardens|
|6iYb2HFDywm3zjuRg...| Restaurants|
|6iYb2HFDywm3zjuRg...|      Bars|
+-----+-----+
only showing top 5 rows
```

## Total Unique Categories

Finally, we are ready to answer the question: **what is the total number of unique categories available?**

Below, implement the code necessary to calculate this figure.

```
In [11]: business_assoc_table.select('categories').distinct().count()
```

```
1330
```

## Top Categories By Business

Now let's find the top categories in this dataset by rolling up categories.

## Counts of Businesses / Category

So now, let's unroll our distinct count a bit and display the per count value of businesses per category.

The expected output should be:

category	count
a	15
b	2
c	45

Or something to that effect.

```
In [12]: business_assoc_table.groupby('categories').count().show(20)
```

```
+-----+-----+
|      categories|count|
+-----+-----+
|   Dermatologists|   351|
| Paddleboarding |    67|
|   Aerial Tours  |     8|
|   Hobby Shops   |   610|
|   Bubble Tea    |   779|
|     Embassy     |     9|
|     Tanning     |   701|
|    Handyman     |   507|
| Aerial Fitness  |    13|
|     Falafel     |   141|
|   Summer Camps  |   308|
|   Outlet Stores |   184|
| Clothing Rental |    37|
| Sporting Goods  |  1864|
| Cooking Schools |   114|
| College Counseling|   20|
| Lactation Services|   47|
| Ski & Snowboard S...|   55|
|      Museums    |   336|
|      Doulas     |    52|
+-----+-----+
only showing top 20 rows
```

## Bar Chart of Top Categories

With this data available, let us now build a barchart of the top 20 categories.

**HINT:** don't forget about the matplotlib magic!

```
%matplotlib plt
```

```
In [13]: top_assoc_business_table = business_assoc_table.groupby('categories').count().orderBy('count', ascending = False)
```

```
In [14]: top_assoc_business_table.show(20)
```

```
+-----+-----+
|      categories|count|
+-----+-----+
```

Restaurants	50763
Food	29469
Shopping	26205
Beauty & Spas	16574
Home Services	16465
Health & Medical	15102
Local Services	12192
Nightlife	11990
Bars	10741
Automotive	10119
Event Planning & ...	9644
Active Life	9231
Coffee & Tea	7725
Sandwiches	7272
Fashion	6599
American (Traditi...	6541
Hair Salons	5900
Pizza	5756
Hotels & Travel	5703
Breakfast & Brunch	5505

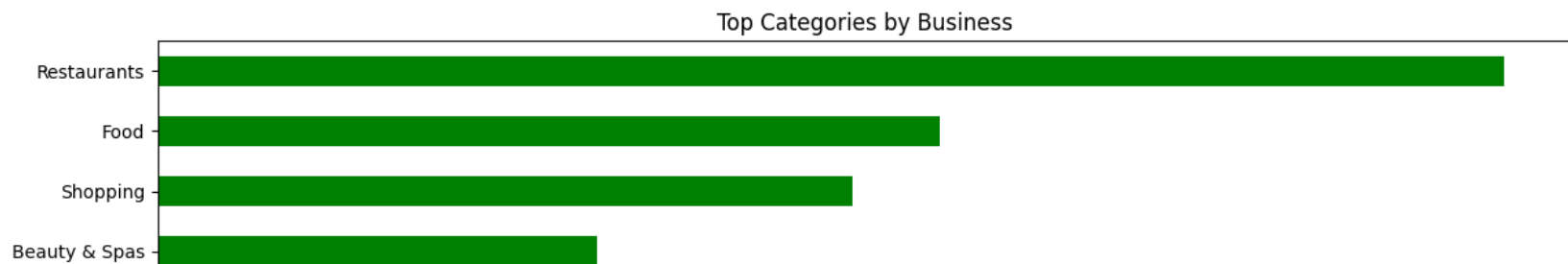
+-----+  
only showing top 20 rows

```
In [15]: top_biz = top_assoc_business_table.limit(20).toPandas()
```

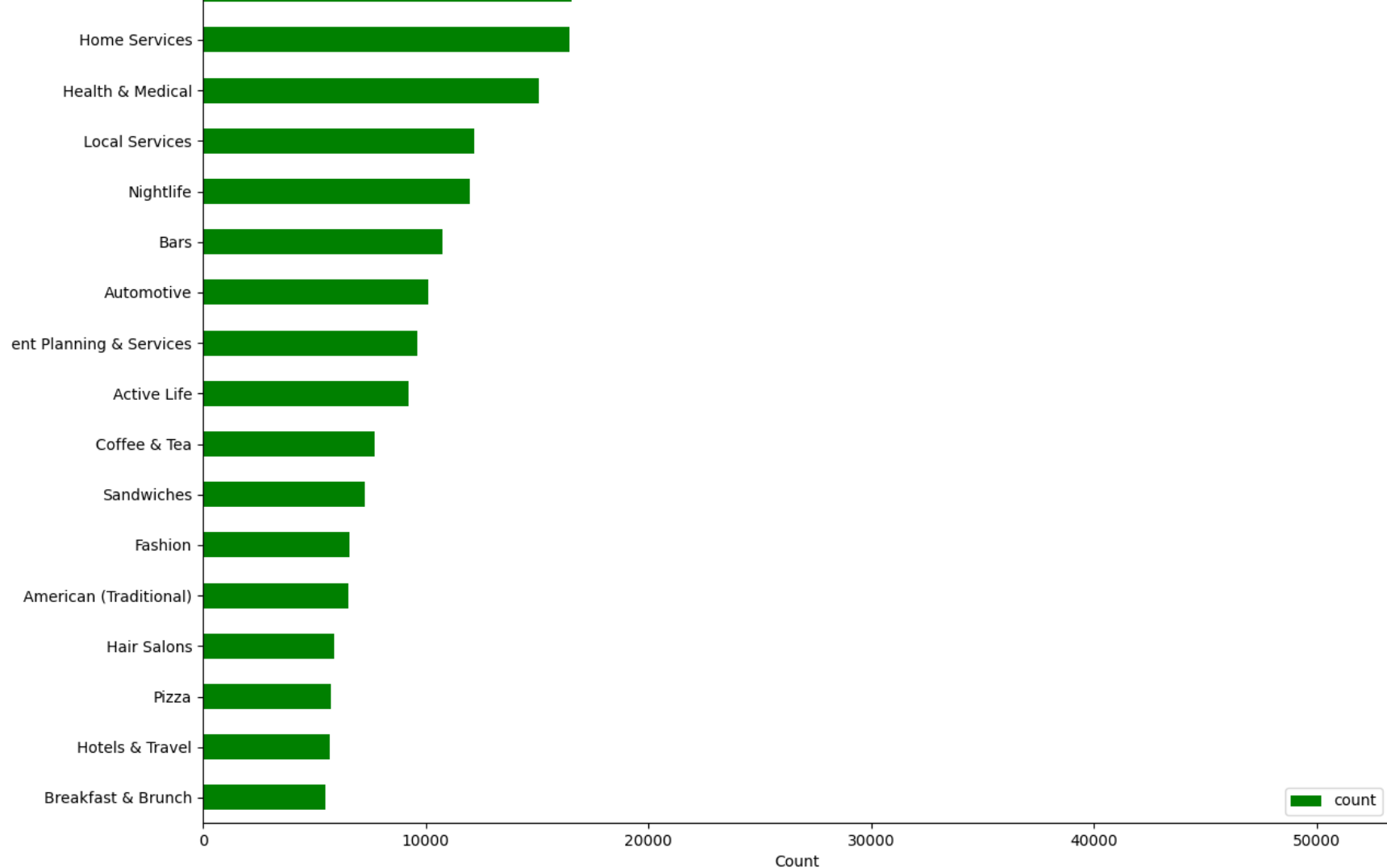
```
In [16]: plt.figure(figsize=(14,12))
top_biz.sort_values(by='count').plot(kind='barh', x='categories', figsize=(14,12), color = 'g')
plt.title('Top Categories by Business')
plt.xlabel('Count')
plt.ylabel('Categories')
```

```
Text(0, 0.5, 'Categories')
```

```
In [17]: %matplotlib plt
```







## Do Yelp Reviews Skew Negative?

Oftentimes, it is said that the only people who write a written review are those who are extremely *dissatisfied* or extremely *satisfied* with the service received.

How true is this really? Let's try and answer this question.

## Loading User Data

Begin by loading the user data set from S3 and printing schema to determine what data is available.

```
In [18]: review_df = spark.read.json('s3://jm-p2-bucket/yelp_academic_dataset_review.json')
review_df.printSchema()
```

```
root
|-- business_id: string (nullable = true)
|-- cool: long (nullable = true)
|-- date: string (nullable = true)
|-- funny: long (nullable = true)
|-- review_id: string (nullable = true)
|-- stars: double (nullable = true)
|-- text: string (nullable = true)
|-- useful: long (nullable = true)
|-- user_id: string (nullable = true)
```

Let's begin by listing the `business_id` and `stars` columns together for the user reviews data.

```
In [19]: review_df.select('business_id', 'stars').show(5)
```

```
+-----+-----+
|      business_id|stars|
+-----+-----+
|buF9druCkbuXLX526...| 4.0|
|RA4V8pr014UyUbDvI...| 4.0|
|_sS2LBIGNT5NQb6PD...| 5.0|
|0AzLzHfOJgLR0whd...| 2.0|
|8zehGz9jnxPqXtOc7...| 4.0|
+-----+-----+
only showing top 5 rows
```

Now, let's aggregate along the `stars` column to get a resultant dataframe that displays *average stars* per business as accumulated by users who **took the time to submit a written review**.

```
In [20]: avg_star_df = review_df.groupby('business_id').avg('stars')
avg_star_df.show(5)
```

```
+-----+-----+
|      business_id|avg(stars)|
+-----+-----+
|uEUweopM30lHcVxj0...|      3.0|
|wdBrDCbZopowEkIEX...|4.538461538461538|
|L3WCfeVozu5etMhz4...|      4.2|
|bOnsvrz1VkbrZM1jV...|      3.8|
|R0IJhEI-zSJpYT1YN...|3.606060606060606|
+-----+-----+
only showing top 5 rows
```

Now the fun part - let's join our two dataframes (reviews and business data) by `business_id`.

```
In [21]: business_review_table = avg_star_df.join(business_df, 'business_id')
```

Let's see a few of these:

```
In [22]: business_review_table.select('avg(stars)', 'stars', 'name', 'city', 'state').show(5)
```

```
+-----+-----+-----+-----+-----+
|      avg(stars)|stars|           name|      city|state|
+-----+-----+-----+-----+-----+
|           5.0|  5.0|   CheraBella Salon|   Peabody|  MA|
|          3.875|  4.0| Mezcal Cantina & ...| Columbus|  OH|
|3.8666666666666667| 4.0|   Red Table Coffee|   Austin|  TX|
|           5.0|  5.0|      WonderWell|   Austin|  TX|
|          3.375|  3.5|     Avalon Oaks|Wilmington|  MA|
+-----+-----+-----+-----+-----+
```

only showing top 5 rows

Compute a new dataframe that calculates what we will call the *skew* (for lack of a better word) between the avg stars accumulated from written reviews and the *actual* star rating of a business (ie: the average of stars given by reviewers who wrote an actual review **and** reviewers who just provided a star rating).

The formula you can use is something like:

$$(\text{row['avg(stars)']} - \text{row['stars']}) / \text{row['stars']}$$

If the **skew** is negative, we can interpret that to be: reviewers who left a written response were more dissatisfied than normal. If **skew** is positive, we can interpret that to be: reviewers who left a written response were more satisfied than normal.

```
In [23]: skewness = business_review_table.select(((business_review_table['avg(stars)'] - business_review_table['stars'])
                                                /business_review_table['stars']))
skewness.show(10)
```

```
+-----+
|((avg(stars) - stars) / stars)|
+-----+
|      0.047619047619047644|
|      0.0303030303030276|
|              0.0|
|      0.05000000000000044|
|      0.037037037037037104|
|      0.008547008547008517|
|      -0.019999999999999...|
+-----+
```

```

|          -0.03703703703703...|
|          -0.05000000000000...|
|          -0.05357142857142857|
+-----+

```

only showing top 10 rows

And finally, graph it!

```

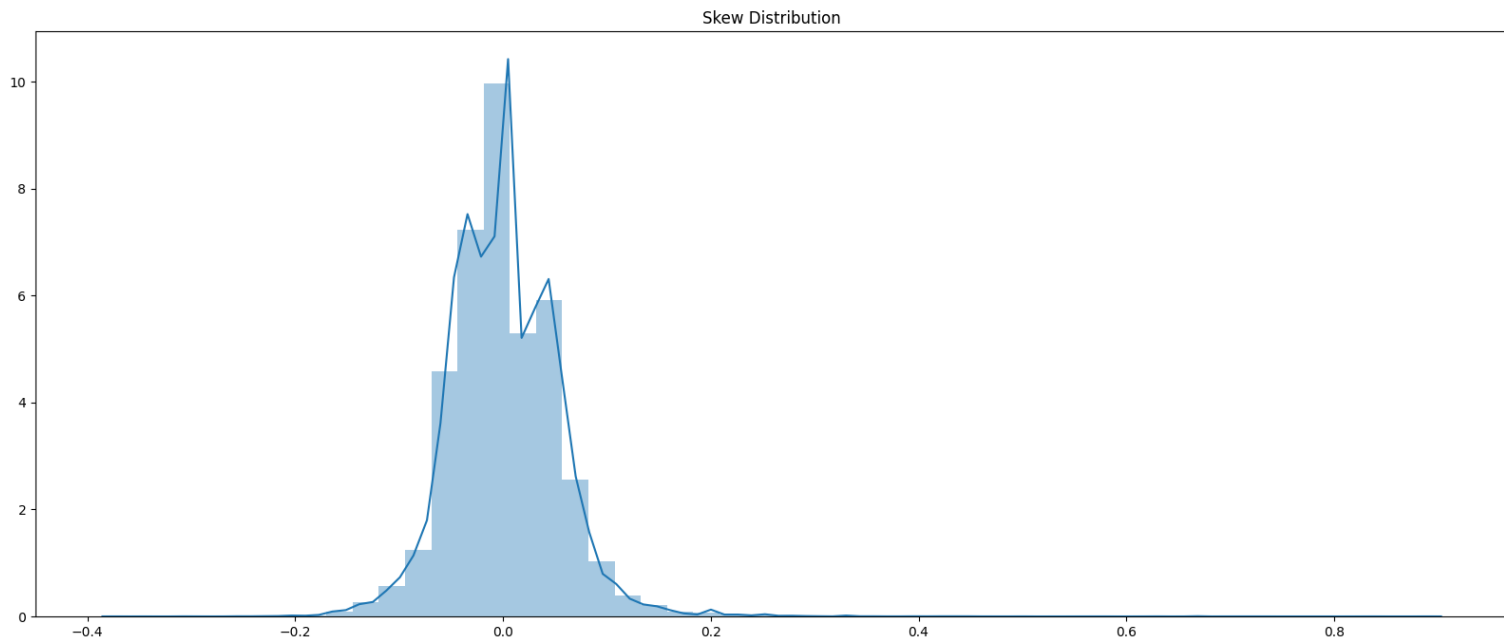
In [24]: skewness_df = skewness.toPandas()
skewness_df.columns=['skewness']

```

```

In [25]: fig, ax = plt.subplots(figsize=(20,8))
skew_plot = sns.distplot(skewness_df)
ax.set_title("Skew Distribution")
%matplotlib plt

```



So, do Yelp (written) Reviews skew negative? Does this analysis actually prove anything? Expound on implications / interpretations of this graph.

Judging from the above plot, the distribution looks to be about normal but leans a little bit more negative. I think in general, people feel a little bit more inclined to leave a review if they have had a negative experience versus if they had a positive experience. If the experience is positive, they will just go back again.

# Should the Elite be Trusted? (Or, some other analysis of your choice)

For the final portion - you have a choice:

- Try and analyze some interesting dimension to this data. The **ONLY** requirement is that you must use the **Users** dataset and join on either the **business\*** or **reviews\*\*** dataset
- Or, you may try and answer the question posed: how accurate or close are the ratings of an "elite" user (check Users table schema) vs the actual business rating.

Feel free to use any and all methodologies at your disposal - only requirement is you must render one visualization in your analysis

## Loading the user json table to my EMR

```
In [26]: user_df = spark.read.json('s3://jm-p2-bucket/yelp_academic_dataset_user.json')
```

```
In [27]: user_df.printSchema()
```

```
root
|-- average_stars: double (nullable = true)
|-- compliment_cool: long (nullable = true)
|-- compliment_cute: long (nullable = true)
|-- compliment_funny: long (nullable = true)
|-- compliment_hot: long (nullable = true)
|-- compliment_list: long (nullable = true)
|-- compliment_more: long (nullable = true)
|-- compliment_note: long (nullable = true)
|-- compliment_photos: long (nullable = true)
|-- compliment_plain: long (nullable = true)
|-- compliment_profile: long (nullable = true)
|-- compliment_writer: long (nullable = true)
|-- cool: long (nullable = true)
|-- elite: string (nullable = true)
|-- fans: long (nullable = true)
|-- friends: string (nullable = true)
|-- funny: long (nullable = true)
|-- name: string (nullable = true)
|-- review_count: long (nullable = true)
|-- useful: long (nullable = true)
|-- user_id: string (nullable = true)
|-- yelping_since: string (nullable = true)
```

In the three below cells I build some dataframes that join the user, review, and business tables together. I also group the elite\_review\_stars together so that I can use this in my skew calculation.

```
In [28]: user_df_clean = user_df.select('user_id', 'elite', 'average_stars', 'review_count')
```

```
In [29]: df_review = review_df.select('user_id', 'business_id', 'stars')
review_elite_df = df_review.join(elite_user_df, 'user_id')
review_elite_df = review_elite_df.withColumnRenamed('stars', 'elite_review_stars')
review_elite_df = review_elite_df.groupby('business_id').avg('elite_review_stars')
```

business_id	actual_business_stars	avg(elite_review_stars)
wdBrDCbZopowEkIEX...	4.5	4.0
L3WCfeVoZu5etMhz4...	4.0	4.6
MPzc6QuEjwk3E3jVT...	3.5	4.3
b0nsvrz1VkbZM1jV...	4.0	5.0
O_BAT_rvszHYBNEM6...	2.5	2.8

only showing top 5 rows

```
In [31]: elite_skewness = business_elite_review_df.select(
    (business_elite_review_df['avg(elite_review_stars)'] - business_elite_review_df['actual_business_stars'])
    /business_elite_review_df['actual_business_stars'])
```

((avg(elite_review_stars) - actual_business_stars) / actual_business_stars)	
	0.14999999999999999
	0.11111111111111111
	-0.11111111111111111
	0.11999999999999993
	0.0
	-0.02832244008714...
	0.06938775510204083
	0.16666666666666666

```
0.08043217286914763|
-0.05263157894736844|
```

```
+-----+
only showing top 10 rows
```

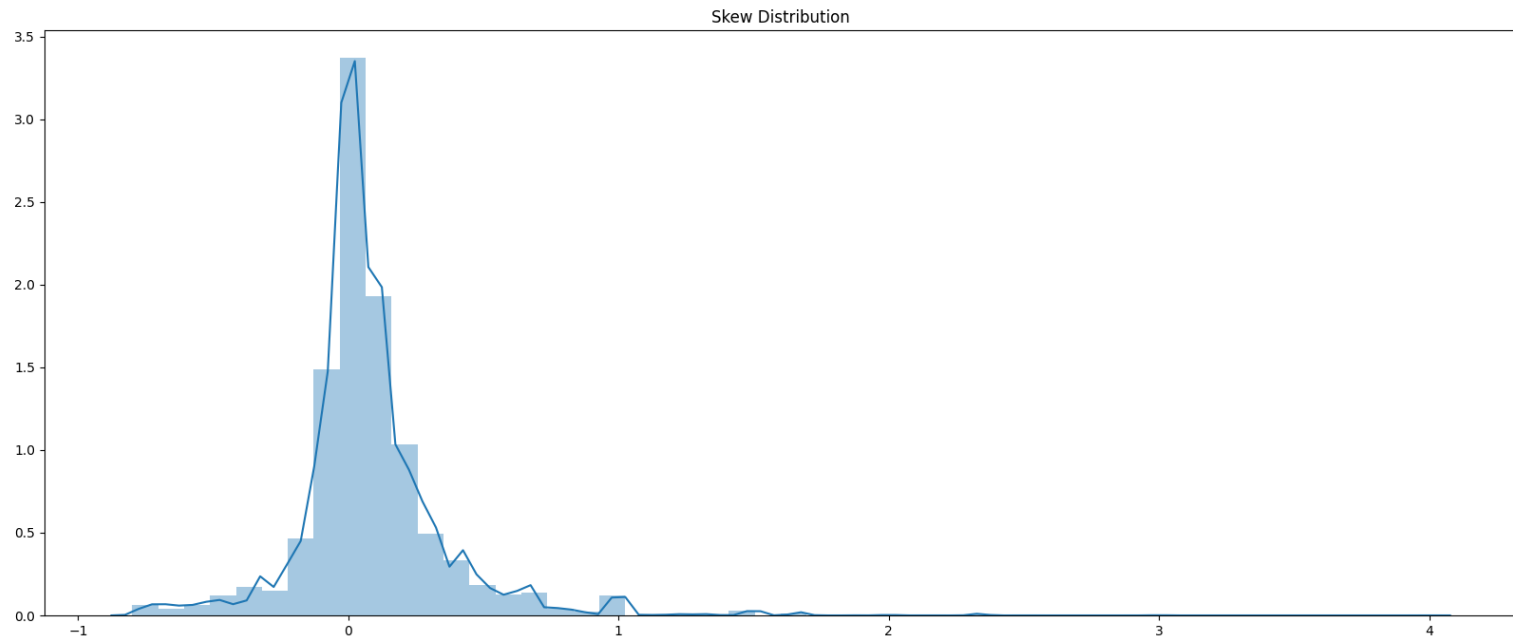
```
In [33]: elite_skewness_df = elite_skewness.toPandas()
         elite_skewness_df.columns=['elite_skew']
```

```
In [34]: elite_skewness_df.head()
```

```
   elite_skew
0   -0.111111
1    0.150000
2    0.228571
3    0.250000
4    0.120000
```

## Plotting the skew

```
In [35]: fig, ax = plt.subplots(figsize=(20,8))
         skew_plot = sns.distplot(elite_skewness_df)
         ax.set_title("Skew Distribution")
         %matplotlib plt
```



## Findings

From the looks of this plot the distribution seems to be about normal with maybe a slightly more negative leaning skew. My judgement is that elite users will have both positive and negative reviews for businesses. I think this is an interesting plot to compare to the plot we did in the previous part. I think elite users will feel more inclined to leave a review whether or not they had a positive or negative experience because they hold this title within Yelp. Their presence in the community makes them feel that reviews in both directions will be useful compared to regular users that may only leave a review if the experience is bad.