

Analyzing the Impact of Compliments on Yelp User Engagement

Problem Statement:

The problem is that we need to understand user behavior on Yelp and the impact of various compliments given to them by other users. Since we are not able to understand user behavior and preferences, it can be challenging to improve the platform experience for both businesses and customers. Not knowing the impact of compliments on users can lead to ineffective marketing strategies and miss opportunities that may come along that way when engaging with customers. The affected parties will include Yelp, business partners, and users. A good starting point would be to analyze user data to answer the questions being asked.

OSEMN Process:

Obtain

The Yelp dataset was obtained from a .tar file, and a python script called "pythonscriptjsoncsv.py" was used to extract it into .json files and a .csv file. The data was then uploaded to a Google Cloud Storage (GCS) bucket named "cs512_group2".

Scrub

Remove NA values from the dataset and other values that may affect the data. Partition the dataset into 160MB size each for dataprep to process successfully.

Explore

Use Dataprep to create a recipe that answers a specific question of interest. In this case, the question is about finding the ratio of "compliment_hot" and "review_count" for Yelp users and sorting them in descending order. I also use BigQuery to partition the dataset to answer the second question. The third question is done in R Studio to create a visualization.

Model

Create a Big Query dataset and a table to calculate the correlation between "compliment_funny" and "compliment_plain". Use a query to count the number of rows and find the correlation coefficient.

Interpret

Use R Studio to visualize the average number of "compliment_cute" being received by users who have been on Yelp before 2010 vs. after 2010. Filter out the dates for those who have been with Yelp before 2010 and after 2010, use the mean function, and remove any NA values. The resulting bar chart shows that those who join Yelp before 2010 have much more compliment_cute than those who joined after.

Description of Process:

Question 1: Which users have the highest ratio of "compliment_hot" to "review_count"?

To start working with the Yelp dataset, I used a Python script called "pythonscriptjsoncsv.py" to extract it from a .tar file into JSON and CSV files. Next, I uploaded the JSON files to my Google Cloud Storage (GCS) bucket, which I named "cs512_group2".

To work with the data in Dataprep, I created a new flow and selected "cloud storage" as the dataset. From there, I chose the "yelp_academic_data_set_user.json" file from my "cs512_group2" bucket. After processing the dataset, I split it into four 1GB chunks for easier handling.

Unfortunately, when attempting to upload the chunks back to the "cs512_group2" bucket, I encountered an issue: the 1GB chunks were too large for Dataprep to handle. To address this, I modified the script to partition the data into smaller chunks (~500MB each) and tried running it through Dataprep again.

Despite these efforts, Dataprep still failed to process the data. As a result, I turned to using the Google Cloud SDK Shell, which prompted me about billing issues.

```
C:\Users\hiimr\Documents\data>gsutil cp yelp_academic_dataset_business.json gs://my-bucket/  
Copying file://yelp_academic_dataset_business.json [Content-Type=application/json]...  
ResumableUploadAbortException: 403 The billing account for the owning project is disabled in state closed
```

Finally, I chose to use PySpark to partition the dataset. I started by loading up PySpark and using the `spark.read.json` method to read in the "yelp_academic_dataset_user.json" file. Through experimentation, I discovered that partitioning the dataset into 160MB files was the largest size that could be processed by Dataprep without crashing. To achieve this, I used PySpark to repartition the "yelp_academic_dataset_user.json" file into 20 parts for better parallelism. Next, I used the `df.write.json` command to write the partitioned data into separate files, resulting in the creation of a new file called "yelp_academic_dataset_user" in the "cs512_group2" bucket. This file contains multiple chunks, each ranging around 160MB in size. The entire process involved multiple stages to complete.

```
>>> df = spark.read.json("gs://cs512_group2/yelp_academic_dataset_user.json")  
[Stage 1:>
```

```
>>> df = df.repartition(20)  
>>> df.write.json("gs://cs512_group2/output/yelp_academic_dataset_user_20_parts.json")
```

After successfully repartitioning the dataset into 20 chunks, I decided to check if the data is appropriately parsed. I went with the “show()” method to display a sample of data. The data should theoretically be parsed and loaded into dictionaries with the schema showing each of the columns as a struct of multiple fields.

```
>>> spark = SparkSession.builder.appName("YelpData").getOrCreate()
>>> yelp_data = spark.read.json("gs://cs512_group2/yelp_academic_dataset_user.json")
[Stage 5:>
[Stage 5:==>
```

I decided to use PySpark's dataframe API to check if the data is appropriately parsed. I load the JSON data from GCS bucket and have it display the first 5 rows of the DataFrame. Since the data is appropriately parsed, I am able to see a tabular display of the first 5 rows for the Yelp data. If it was not appropriately parsed, then I may see errors or weird outputs when displaying the data.

```
>>> yelp_data.show(5)
[Stage 6:>

+-----+-----+-----+-----+
|-----+-----+-----+-----+
|-----+-----+-----+-----+
|-----+-----+-----+-----+
|--+---+---+---+---+---+---+---+---+
|----+---+---+---+---+---+---+---+
|-----+-----+-----+-----+

|average_stars|compliment_cool|compliment_cute|co
mpliment_funny|compliment_hot|compliment_list|com
pliment_more|compliment_note|compliment_photos|co
mpliment_plain|compliment_profile|compliment_writ
er| cool| elite|fans| f
riends|funny| name|review_count|useful|
      user id|       yelping since|
```

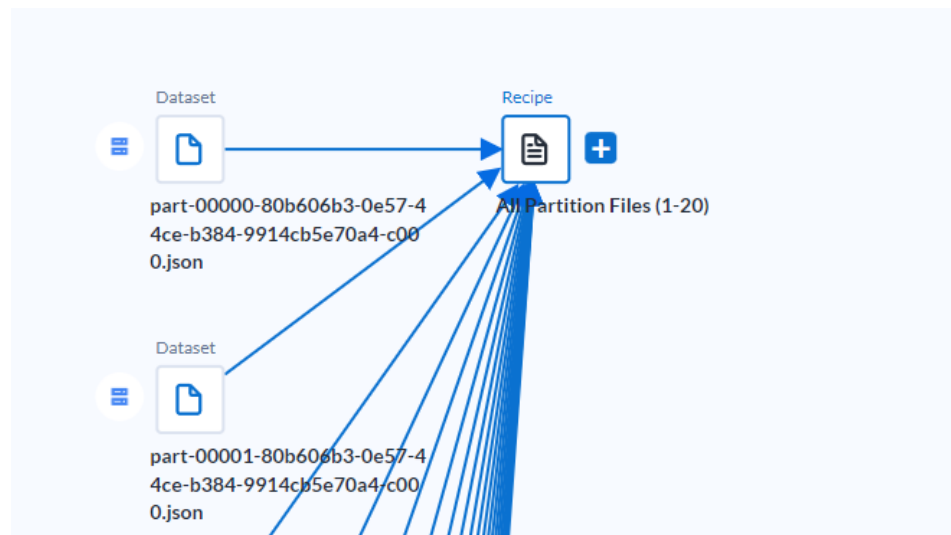
The final step is to edit the recipe in dataprep after successfully uploading all of the partitions.

The screenshot shows the Dataprep interface. On the left, a file browser shows a directory path `/yelp_academic_dataset_user/part-` with a search bar and a "Create Dataset with Parameters" button. Below this is a table of uploaded partitions:

NAME	SIZE	LAST UPDATED
+ [icon] part-00000-80b606b3-0...	162.12MB	Today at 9:59 PM
+ [icon] part-00001-80b606b3-0...	159.89MB	Today at 9:59 PM
+ [icon] part-00002-80b606b3-0...	161.18MB	Today at 9:59 PM
+ [icon] part-00003-80b606b3-0...	160.22MB	Today at 9:59 PM
+ [icon] part-00004-80b606b3-0...	160.11MB	Today at 9:59 PM
+ [icon] part-00005-80b606b3-0...	159.45MB	Today at 9:59 PM

On the right, a panel titled "20 New Datasets" shows a preview of a dataset with columns `compliment_cool` and `compliment_hot`. Below this, a specific dataset is selected: `part-00002-80b606b3-0e57-44ce-b384-9914cb5e70a4-c000.json`.

I accessed the "Final Project Flow" view and selected all of the partitions. I then used the "union" option to combine all of the .json files into a single recipe. Using this recipe, I was able to answer my question of interest. During the recipe editing process, I could preview my newly created column called "hot_review_ratio", which calculates the ratio of "compliment_hot" to "review_count".



Final Project > All Partition Files (1-20)
Job 18360643
Finished Today at 11:47 PM

[Overview](#) [Output destinations](#) [Profile](#) [Dependency graph](#) [Data sources](#)

Output data

A_c^B	column2	1_3^2	column3	1_3^2	column4	1_3^2	column5
ErEfjm4wKBV0bfIFBvva6Q	0			1		0	
zB3K0dWxbNVQuuEqQnW4g	0			1		0	
WcFACH4J--04e0dXc3a7g	0			57		0	
W0sU93vc0G0GehAqVZuusQ	0			2		0	
nDH0aEV0ki5FYfUERn_rEw	1			78		0.0128205128205128	
XnFMFCm1V0kBO_pSQLVuzg	0			5		0	
A0pWKA2NYgN-1J4Nzgco2w	0			9		0	
BXQDsVamsDh44Wh4ygp r9g	0			1		0	
VF_-64Xe8uRfh93pnJ1haw	2			73		0.0273972602739726	
XavKeY7PmNchxvYfCtFhRn	0			0		0	

4 columns 1.988M rows This is a preview of the current data in your destination. It might not reflect the output from this particular job run.


[View on Cloud Storage](#) [Download](#) [View details](#)

Execution stages


✓	Transform with profile Completed Today at 11:47 PM, started Today at 11:32 PM • Ran for 15 min Environment Dataflow <div><div></div><div>98% valid values</div><div>2% mismatching values</div><div>< 1% missing values</div></div> View steps and dependencies View profile View dataflow job
✓	Publish Completed Today at 11:47 PM, started Today at 11:47 PM • Ran for 20 sec Activity <div><div></div>All Partition Files (1-20).csv View all</div>

Question 2: Is there a correlation between the number of "compliment_funny" a user receives and the number of "compliment_plain" they receive? If so, what is the correlation coefficient?

I went into Dataproc to create a cluster, and after successfully creating it, I created a job with a Python script I had prepared. I uploaded the Python script to my bucket and used the Hadoop jar file to connect to BigQuery. Once the job was successfully executed, I proceeded to BigQuery to calculate the correlation between the "compliment_funny" and "compliment_plain" columns.

Job ID	job-abc04cc5
Job UUID	b785e4f5-4319-41ae-9efe-e83ee9c9f631
Type	Dataproc Job
Status	 Succeeded










Output
LINE WRAP: OFF
⌵

 Spark jobs take ~60 seconds to initialize resources.
 DISMISS

```

at com.google.common.util.concurrent.AbstractFuture.get(AbstractFuture.java:516)
at com.google.common.util.concurrent.FluentFuture$TrustedFuture.get(FluentFuture$TrustedFuture.java:55)
at org.apache.hadoop.util.concurrent.ExecutorHelper.logThrowableFromAfterExecute(ExecutorHelper.java:116)
at org.apache.hadoop.util.concurrent.HadoopThreadPoolExecutor.afterExecute(HadoopThreadPoolExecutor.java:116)
at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:111)
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:602)
at java.lang.Thread.run(Thread.java:750)
23/03/15 20:29:12 WARN org.apache.hadoop.util.concurrent.ExecutorHelper: Thread (Thread-1) interrupted
java.lang.InterruptedException
at com.google.common.util.concurrent.AbstractFuture.get(AbstractFuture.java:516)
at com.google.common.util.concurrent.FluentFuture$TrustedFuture.get(FluentFuture$TrustedFuture.java:55)
at org.apache.hadoop.util.concurrent.ExecutorHelper.logThrowableFromAfterExecute(ExecutorHelper.java:116)
at org.apache.hadoop.util.concurrent.HadoopThreadPoolExecutor.afterExecute(HadoopThreadPoolExecutor.java:116)
at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:111)
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:602)
at java.lang.Thread.run(Thread.java:750)
23/03/15 20:29:12 INFO org.apache.hadoop.mapred.FileInputFormat: Total input files: 1
23/03/15 20:30:31 INFO com.google.cloud.hadoop.repackaged.gcs.com.google.cloud.hadoop.fs.gcs.GoogleHadoopFileSystem: Stopped SpillWriter
23/03/15 20:30:31 INFO org.sparkproject.jetty.server.AbstractConnector: Stopped SpillWriter

```

Buckets > cs512_group2 > user_chunks > user 			
UPLOAD FILES		UPLOAD FOLDER	CREATE FOLDER
TRANSFER DATA		MANAGE HOLDS	DELETE
Filter by name prefix only ▼		 Filter	Filter objects and folders  Show deleted data
<input type="checkbox"/>	Name	Size	Type
<input type="checkbox"/>	 _SUCCESS	0 B	application/octet-
<input type="checkbox"/>	 part-00000	159.2 MB	application/octet-
<input type="checkbox"/>	 part-00001	162.1 MB	application/octet-
<input type="checkbox"/>	 part-00002	158.9 MB	application/octet-
<input type="checkbox"/>	 part-00003	160.1 MB	application/octet-
<input type="checkbox"/>	 part-00004	159 MB	application/octet-

To answer the second question, I chose to utilize Big Query. First, I created a new dataset and table, uploading the “user_chunks” data from my Google Cloud Storage bucket. I then constructed a query to count the number of rows and apply the correlate function to determine the correlation between the number of "compliment_funny" a user receives and the number of "compliment_plain". The resulting correlation coefficient was 0.867, indicating a strong positive correlation. This suggests that users who receive many "compliment_funny" also tend to receive many "compliment_plain", and vice versa.

<pre> 1 SELECT 2 COUNT(*) OVER() AS row_count, 3 CORR(compliment_funny, compliment_plain) AS correlation_coefficient 4 FROM 5 dataset_table_1.yelp_table_json 6 WHERE 7 compliment_funny IS NOT NULL 8 AND compliment_plain IS NOT NULL; </pre>			
Query results			
JOB INFORMATION		RESULTS	JSON EXECUTION DETAILS EXECUTION GRAPH PREVIEW
Row	row_count	correlation_coef	
1	1	0.86710332...	

Question 3: What is the average number of "compliment_cute" received by users who have been on Yelp since before 2010, and how does this compare to the average number of "compliment_cute" received by users who joined Yelp after 2010?

For the last question, I used R Studio to visualize the average number of "compliment_cute" received by Yelp users who joined before 2010 and those who joined after. I first filtered the data for users who joined before 2010 and those who joined after, then used the mean function to calculate the average number of "compliment_cute" received. I removed any NA values that could affect the data. I created a bar chart to display the results, and while I tried to improve its appearance, this is the final version. The chart shows that users who joined Yelp before 2010 received significantly more "compliment_cute" than those who joined after.

```
# Filter the data for users before 2010 and after 2010 timeframe
users_before_2010 <- yelp_data |>
  filter(as.Date(yelping_since) < as.Date("2010-01-01"))
users_after_2010 <- yelp_data |>
  filter(as.Date(yelping_since) >= as.Date("2010-01-01"))

# Find the average number of compliment_cute for each group
avg_cute_before_2010 <- mean(users_before_2010$compliment_cute, na.rm = TRUE)
avg_cute_after_2010 <- mean(users_after_2010$compliment_cute, na.rm = TRUE)

cat("Average compliment_cute before 2010:", avg_cute_before_2010, "\n")
cat("Average compliment_cute after 2010:", avg_cute_after_2010, "\n")

# Visualization to compare the averages
df <- data.frame(
  group = c("Before 2010", "After 2010"),
  avg_cute = c(avg_cute_before_2010, avg_cute_after_2010)
)

ggplot(df, aes(x = group, y = avg_cute)) +
  geom_bar(stat = "identity", fill = "#FF5733", width = 0.6) +
  geom_text(aes(label = round(avg_cute, 2)), vjust = -0.5, color = "white", size = 8) +
  theme_minimal() +
  theme(
    plot.title = element_text(face = "bold", size = 20, hjust = 0.5),
    axis.text = element_text(size = 15),
    axis.title = element_text(face = "bold", size = 16),
    legend.title = element_blank(),
    legend.position = "none"
  ) +
  labs(title = "Average Compliment_cute by Yelp Join Date",
       x = "Join Date", y = "Average Compliment_cute")
```

