

Project Report: Yelp Dataset

Group 2.

Adriana Marcela Penaranda Baron

Carlos Rey Pinto

Eduardo Roberto Williams Cascante

Haldo Jose Somoza Solis

Ishika Sukhija

Luis Alejandro Gutierrez Hayek

Marzieh Mohammadi Kokaneh

Nilesh Khurana

Lambton College

Big Data Analytics

BDM 1024 - Data Technology Solutions

To: Mr. Bhavik Gandh

Mississauga, Ontario

16 Aug, 2023

Table of Contents

Data Model Diagram.....	3
Entities	3
Decision Tree Analysis.....	5
Entity relationship diagram (ER Diagram).....	7
Yelp Dataset Dictionary.....	8
Yelp Dataset Analysis in a PostgreSQL Database	12
Python program created to transform the information from json format to PostgreSQL.....	12
PostgreSQL database and tables with relationships created in the Cloud (Azure).	13
Dataset (json documents) imported to the Azure database.	14
Dataset cleaning process	15
Dataset analysis and insights	18
Map Reduce Part	42
Spark (PySpark) Part.....	29
Conclusions and Recommendations	65
References	69

Project Report: Yelp Dataset

Data Model Diagram

The dataset selected is the "Yelp Open Dataset," available as JSON files on the Yelp website and with a size of 8.65GB uncompressed. (yelp, 2023)

Entities

First, the following entities were identified in the selected Yelp company dataset:

- Businesses: Contains information related to companies, like name, city, address, etc.
- Attributes: Contains different features related to businesses.
- Categories: Contains the types of businesses
- Hours: Contains schedules, including days of the week and hours that operate the business.
- Check-ins: Enumerates check-ins for a business.
- Photos: This page contains some pictures of the businesses.
- Users: This contains the information of the customers that review or leave tips to improve the business service.
- Users Elite: This contains the number of years that the user has had elite status.
- Friends: Contains the relationships between users.
- Tips: Contains a short recommendation that users give to businesses for them to improve their services.
- Reviews: Contains full review text data, including the user that wrote the review and the business the review is written for.

The following table describes the relationships between “Businesses” entities:

Entity	Relationship type	Cardinality
Attributes	Composite association regarding Businesses	Businesses → Attributes: One to many
Categories	Composite association regarding Businesses	Businesses → Categories: One to many
Hours	Composite association regarding Businesses	Businesses → Hours: One to many
Checkins	Composite association regarding Businesses	Businesses → Checkins: One to many
Photos	Composite association regarding Businesses	Businesses → Photos: One to many

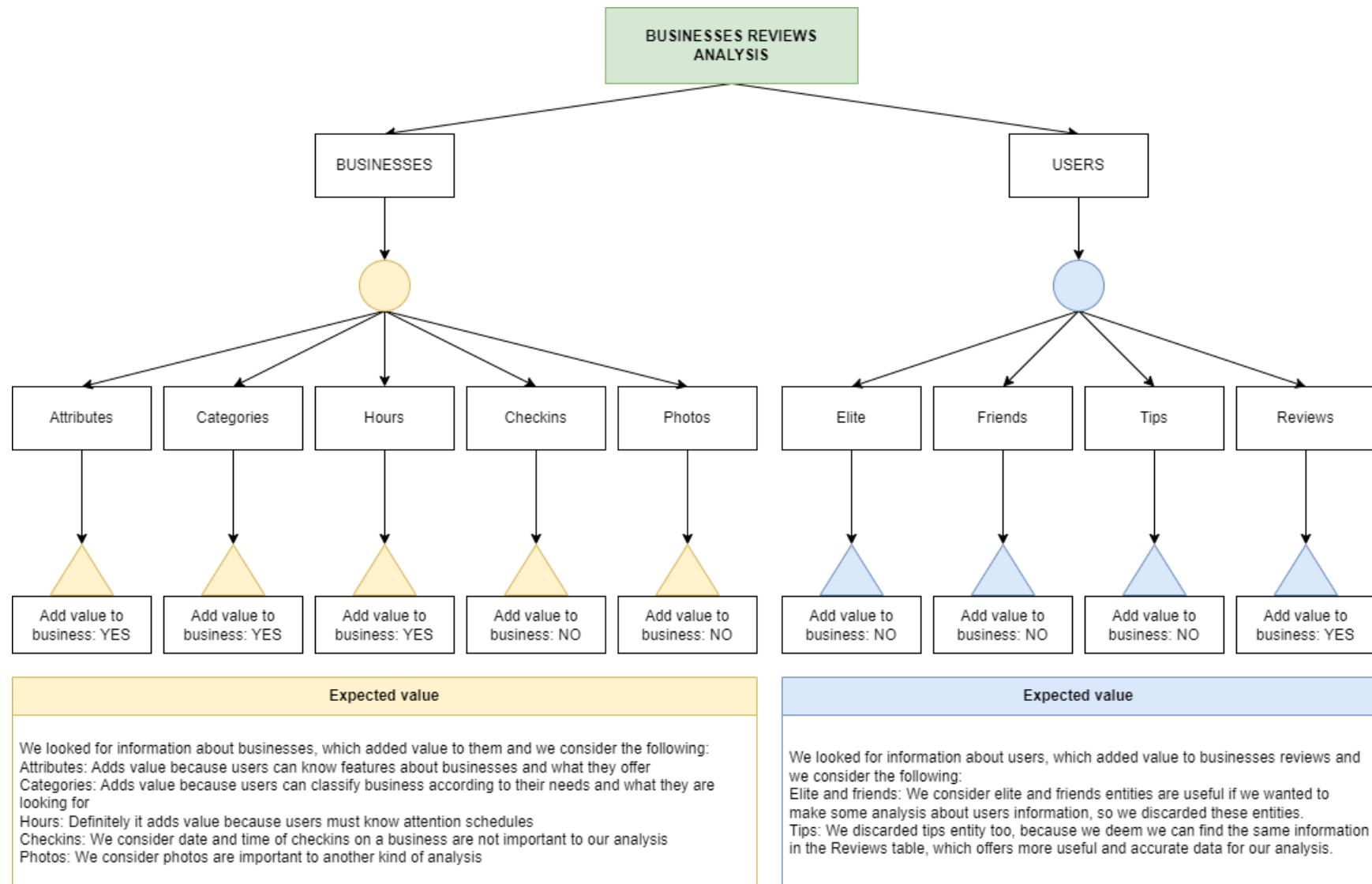
The following table describes the relationships between “Users” entities:

Entity	Relationship type	Cardinality
Users “elite”	Composite association regarding Users	Users → Users elite: One to many
Friends	Composite association regarding Users	Users → Friends: One to many

In addition, it was identified an aggregation association between “Businesses” and “Users” entities, which is many to many. This relationship is represented through Tips and Reviews entities, as follow:

Entity	Relationship type	Cardinality
Tips	Composite association regarding Users and Businesses	Users → Tips: One to many Businesses → Tips: One to many
Reviews	Composite association regarding Users and Businesses	Users → Reviews: One to many Businesses → Reviews: One to many

Decision Tree Analysis



Following a qualitative decision tree analysis, it was agreed to work with the following 6 entities, which contain the main features of the businesses and meet all requirements for SQL database importation:

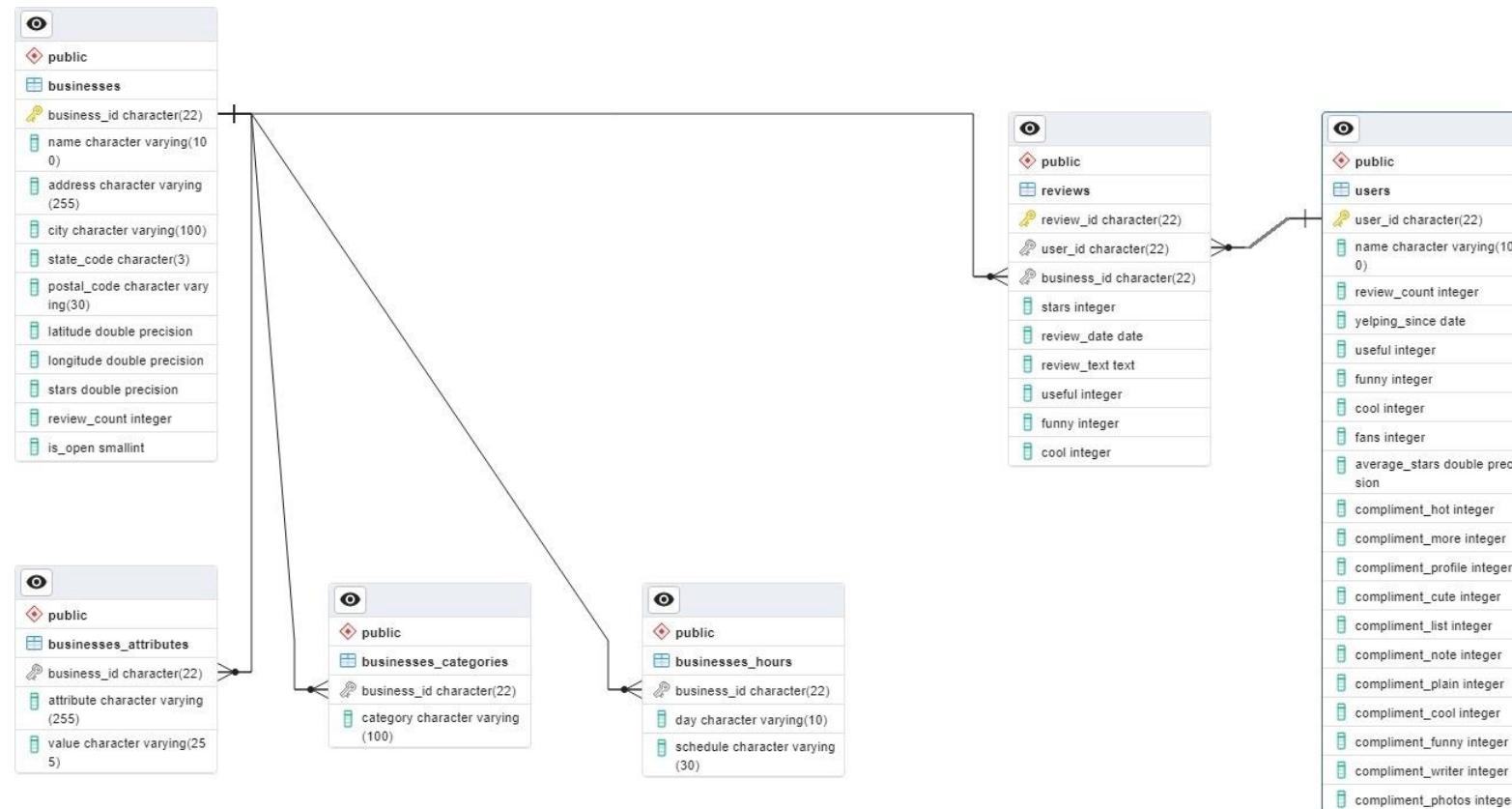
- Businesses
- Attributes
- Categories
- Hours
- Users
- Reviews

The following table contains the main reasons why the entities about business and users were or were not selected:

Entity	Reason
Businesses	The principal entity that contains the main description columns about the business
Attributes	This table adds value to the data because users can learn about the businesses features and what they offer, and it would be possible to analyze the preferred features for users and categorize them.
Categories	It is very important to know how businesses are classified according to what they offer.
Hours	It is important that users know the business hours, and the most demanded business hours can be part of the analysis.
Check-ins	This entity will be considered in the final project to process it in a NoSQL Database
Photos	Photos are very important to the Yelp database, but they do not contribute to this SQL database analysis.
Elite and friends	These entities will probably be considered in the final project to process it in a NoSQL Database based on the type of information it contains.
Tips	

Note: Attributes, integrity rules like PK and FK, and Indices can be found in the data dictionary.

Entity relationship diagram (ER Diagram)



Note: ER Diagram was creating using PgAdming4

Yelp Dataset Dictionary

The dataset dictionary was created. This step is crucial for the creation of the database and the safe import of the data, always maintaining the information integrity. This data dictionary was created considering the Yelp documentation provided by each entity and serves as a valuable reference since it contains details such as the variables, data types, definitions, and some examples. This documentation facilitates the analysis of the data, minimizes misinterpretation of the data since there is a better understanding of the information, and allows more reliable and effective conclusions and recommendations to be reached.

businesses					
Description:	Contains business data including location data, attributes, and categories				
Field name	Data type	Key	Constraint	Description	Example
business_id	char(22)	PK	not null	business id	tnhfDv5Il8EaGSXZGiuQGg
name	varchar(100)		null	the business's name	Garaje
address	varchar(255)		null	the full address of the business	475 3rd St
city	varchar(100)		null	the city	San Francisco
state_code	char(3)		null	state code	CA
postal_code	varchar(30)		null	the postal code	94107
latitude	float		null	latitude	37.7817529521
longitude	float		null	longitude	-122.39612197
stars	float		null	star rating, rounded to half-stars	4.5
review_count	int		null	number of reviews	1198
is_open	smallint		null	0 or 1 for closed or open, respectively	1
Index	businesses_idx0 businesses_idx1 (state_code, city, is_open)				

businesses_attributes					
Description: Attributes and subattributes of businesses					
Field name	Data type	Key	Constraint	Description	Example
business_id	char(22)	FK	not null references businesses (business_id)	business id	tnhfDv5Ii8EaGSXZGiUQGg
attribute	varchar(255)		not null	Attribute:Subattribute description	RestaurantsTakeOut BusinessParking:garage
value	varchar(255)		not null	true / false	
Index	businesses_attributes_idx0 (business_id, checked)				

businesses_categories					
Description: Categories of businesses					
Field name	Data type	Key	Constraint	Description	Example
business_id	char(22)	FK	not null references businesses (business_id)	business id	tnhfDv5Ii8EaGSXZGiUQGg
category	varchar(100)		not null	category description	Mexican Burgers
Index	businesses_categories_idx0 (business_id, category)				

businesses_hours					
Description: Businesses schedules					
Field name	Data type	Key	Constraint	Description	Example
business_id	char(22)	FK	not null references businesses (business_id)	business id	tnhfDv5Ii8EaGSXZGiUQGg
day	varchar(10)		not null check(day in ('Sunday', 'Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday'))	day of week	Tuesday
schedule	varchar(30)		not null	hour from - hour to	10:00-21:00
Index	businesses_hours_idx0 (business_id)				

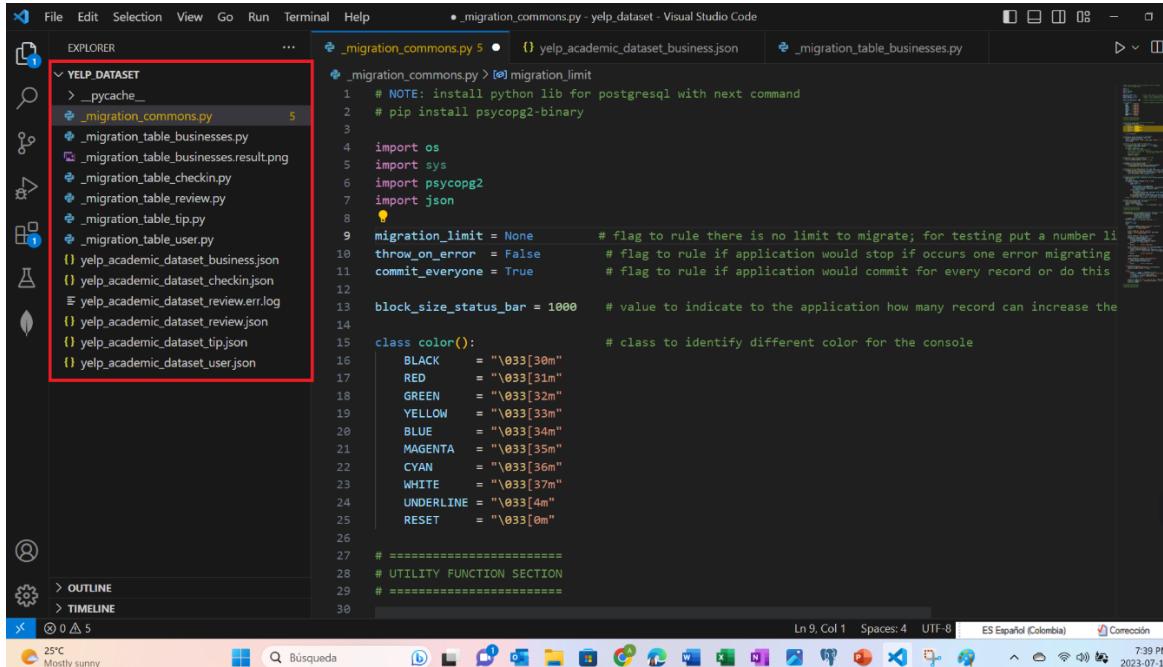
Users					
Description:	User data including all the metadata associated with the user				
Field name	Data type	Key	Constraint	Description	Example
user_id	char(22)	PK	not null	user id	Ha3iJu77CxlrFm-vQRs_8g
name	varchar(100)		null	the user's first name	Sebastien
review_count	int		null	the number of reviews they've written	56
yelping_since	date		null	when the user joined Yelp, formatted like YYYY-MM-DD	2011-01-01
useful	int		null	number of useful votes sent by the user	21
funny	int		null	number of funny votes sent by the user	88
cool	int		null	number of cool votes sent by the user	15
fans	int		null	number of fans the user has	1032
average_stars	float		null	average rating of all reviews	4.31
compliment_hot	int		null	number of hot compliments received by the user	339
compliment_more	int		null	number of more compliments received by the user	668
compliment_profile	int		null	number of profile compliments received by the user	42
compliment_cute	int		null	number of cute compliments received by the user	62
compliment_list	int		null	number of list compliments received by the user	37
compliment_note	int		null	number of note compliments received by the user	356
compliment_plain	int		null	number of plain compliments received by the user	68
compliment_cool	int		null	number of cool compliments received by the user	91
compliment_funny	int		null	number of funny compliments received by the user	99
compliment_writer	int		null	number of writer compliments received by the user	95
compliment_photos	int		null	number of photo compliments received by the user	50
Index	users_idx0 users_idx1 (yelping_since)				(name)

Reviews					
Description: Contains full review text data including the user_id that wrote the review and the business_id the review is written for					
Field name	Data type	Key	Constraint	Description	Example
review_id	char(22)	PK	not null	review id	zdSx_SD6obEhz9VrW9uAWA
user_id	char(22)	FK	not null references users (user_id)	user id	Ha3iJu77CxlrFm-vQRs_8g
business_id	char(22)	FK	not null references businesses (business_id)	business id	tnhfDv5II8EaGSXZGiuQGg
stars	int		null	star rating	4
review_date	date		null	date formatted YYYY-MM-DD	2016-03-09
review_text	text		null	the review itself	Great place to hang out after work: the prices are decent, and the ambience is fun. It's a bit loud, but very lively. The staff is friendly, and the food is good. They have a good selection of drinks
useful	int		null	number of useful votes received	0
funny	int		null	number of funny votes received	0
cool	int		null	number of cool votes received	0
Index	reviews_idx0 reviews_idx1 reviews_idx2 (review_date, business_id, user_id)				(user_id, (business_id, review_date) review_date)

Yelp Dataset Analysis in a PostgreSQL Database

Python program created to transform the information from json format to PostgreSQL

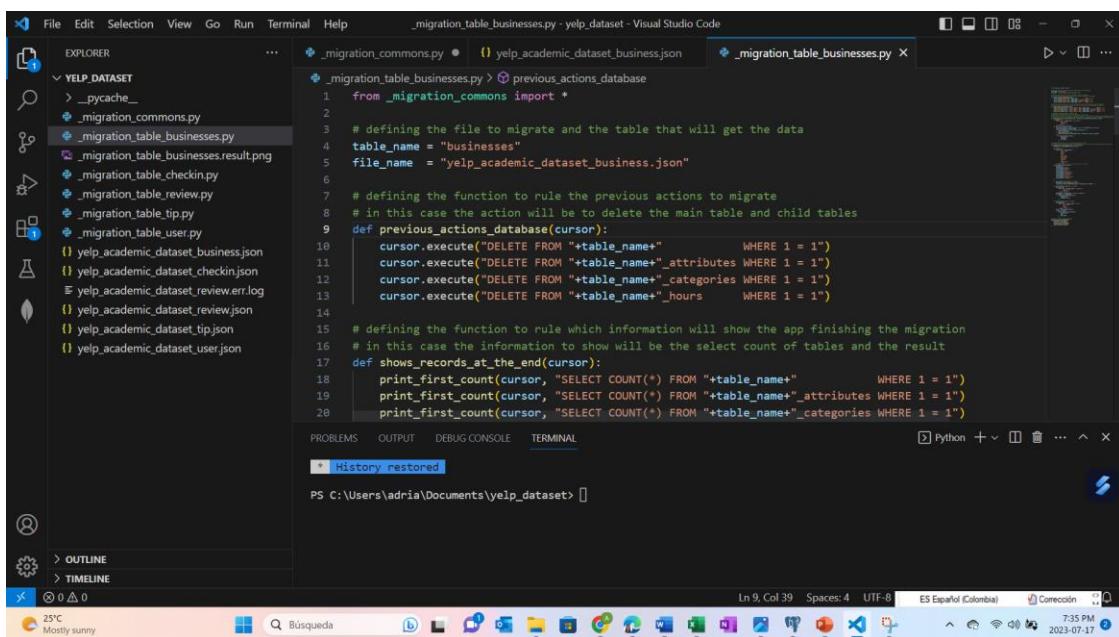
It was necessary to create 7 python programs to transform the dataset from json format to PostgreSQL. The file “_migration_commons.py” contains the main features of the data transformation and connection. Also, it was created one py file for each entity, as shown below:



```

File Edit Selection View Go Run Terminal Help _migration_commons.py - yelp_dataset - Visual Studio Code
EXPLORER _migration_commons.py 5 | yelp_academic_dataset_business.json | _migration_table_businesses.py ...
_migration_commons.py > migration_limit
1 # NOTE: install python lib for postgresql with next command
2 # pip install psycopg2-binary
3
4 import os
5 import sys
6 import psycopg2
7 import json
8
9 migration_limit = None      # flag to rule there is no limit to migrate; for testing put a number like
10 throw_on_error = False     # flag to rule if application would stop if occurs one error migrating
11 commit_everyone = True      # flag to rule if application would commit for every record or do this
12
13 block_size_status_bar = 1000 # value to indicate to the application how many record can increase the
14
15 class color():
16     BLACK    = "\033[30m"
17     RED     = "\033[31m"
18     GREEN   = "\033[32m"
19     YELLOW  = "\033[33m"
20     BLUE    = "\033[34m"
21     MAGENTA = "\033[35m"
22     CYAN   = "\033[36m"
23     WHITE   = "\033[37m"
24     UNDERLINE = "\033[4m"
25     RESET   = "\033[0m"
26
27 # =====
28 # UTILITY FUNCTION SECTION
29 # =====
30

```



```

File Edit Selection View Go Run Terminal Help _migration_table_businesses.py - yelp_dataset - Visual Studio Code
EXPLORER _migration_commons.py > previous_actions_database | _migration_table_businesses.py ...
_migration_table_businesses.py > previous_actions_database
1 from _migration_commons import *
2
3 # defining the file to migrate and the table that will get the data
4 table_name = "businesses"
5 file_name = "yelp_academic_dataset_business.json"
6
7 # defining the function to rule the previous actions to migrate
8 # in this case the action will be to delete the main table and child tables
9 def previous_actions(cursor):
10     cursor.execute("DELETE FROM "+table_name+" WHERE 1 = 1")
11     cursor.execute("DELETE FROM "+table_name+"_attributes WHERE 1 = 1")
12     cursor.execute("DELETE FROM "+table_name+"_categories WHERE 1 = 1")
13     cursor.execute("DELETE FROM "+table_name+"_hours WHERE 1 = 1")
14
15 # defining the function to rule which information will show the app finishing the migration
16 # in this case the information to show will be the select count of tables and the result
17 def show_records_at_the_end(cursor):
18     print_first_count(cursor, "SELECT COUNT(*) FROM "+table_name+" WHERE 1 = 1")
19     print_first_count(cursor, "SELECT COUNT(*) FROM "+table_name+"_attributes WHERE 1 = 1")
20     print_first_count(cursor, "SELECT COUNT(*) FROM "+table_name+"_categories WHERE 1 = 1")

```

PostgreSQL database and tables with relationships created in the Cloud (Azure).

When carrying out this project, the opportunity was taken to put concepts from the classes into practice, such as the use of the Cloud, for which it was decided to create the database using Azure services, where each group participant had access to the same transformed dataset, run queries to understand the dataset and participate in the whole process being able to produce some insights.

The screenshot shows the Microsoft Azure portal with the URL portal.azure.com/#@mylambton.onmicrosoft.com/resource/subs.... The page title is "postgresql-server-01 | Connect". The left sidebar lists various management options like Overview, Activity log, Access control (IAM), Tags, Diagnose and solve problems, Migration, Compute + storage, Networking, Databases, Connect (which is selected), Server parameters, Replication, Maintenance, High availability, Backup and restore, Advisor recommendations, Locks, Power Platform (Power BI (preview)), Security (Data encryption, Authentication), and Power Platform (Power BI (preview)).

The main content area displays a "Pre-requisites check" section with a green checkmark next to "Server is in Ready state". Below it is a "Database name" dropdown. The "Connection details" section is highlighted with a red box. It contains instructions to set environment variables for WSL, Azure Cloud Shell, etc., followed by a code block:

```
export PGHOST=postgresql-server-01.postgres.database.azure.com
export PGUSER=superuser
export PGPORT=5432
export PGDATABASE=
export PGPASSWORD="(your-password)"
```

Below this, another section says: "After setting these variables, you can connect to your database server using various PostgreSQL utilities (psql, pg_dump, pg_restore, pgbench, createdb) without specifying connection options. For example, you can now simply type psql to connect:" followed by a code block:

```
psql
```

Dataset (json documents) imported to the Azure database.

Once the Python program and the database were ready, it was time to import the data to the Azure database through PgAdmin4. The following image shows the migration process for the table users:

PostgreSQL driver for Python installation

```
C:\Windows\system32\cmd.exe
C:\Windows [Versión 10.0.19045.3086]
) Microsoft Corporation. Todos los derechos reservados.

:Users\eduardo pip install psycopg2-binary
Collecting psycopg2-binary
  Downloading psycopg2-binary-2.9.6-cp311-cp311-win_amd64.whl (1.2 MB)
    1.2/1.2 MB 9.2 MB/s eta 0:00:00
Installing collected packages: psycopg2-binary
Successfully installed psycopg2-binary-2.9.6
:Users\eduardo>
```

Running python program migration

Finishing process

table_user.py [migration_commons]

```
from __migration_commons import *
...
1 table_name = "users"
2 file_name = "yelp_academic_dataset_user.json"
3
4 def previous_actions_database(cursor):
5     pass
6     cursor.execute("DELETE FROM "+table_name+" WHERE 1 = 1")
7     cursor.execute("DELETE FROM "+table_name+"_friends WHERE 1 = 1")
8     cursor.execute("DELETE FROM "+table_name+"_elite WHERE 1 = 1")
9
10 def show_records_at_the_end(cursor):
11     print(first_count(cursor, "SELECT COUNT(*) FROM "+table_name+" WHERE 1 = 1"))
12
13 def main():
14     previous_actions_database(cursor)
15     show_records_at_the_end(cursor)
```

Shell

```
Results: OK => Total: 1977957, Errors: 0
```

PAGE 05 showing migrated records

```
    SELECT COUNT(*) FROM users      WHERE 1 = 1 => 1977957 rows
```

Closing db objects (cursor and connection)

result: OK

If errors, check the file: YELP_ACADEMIC_DATASET_USER.ERR.LOG

- - - FINISHING EXECUTION - - -

Users table verification

The screenshot shows the pgAdmin interface with the following details:

- Toolbar:** Dashboard, Properties, SQL, Statistics, Dependencies, Dependents, Processes.
- Search Bar:** yelp_dataset/superuser@dbReviews
- Filter Bar:** No limit, with dropdown and various icons.
- Query Tab:** Active, showing "Query History".
- Query Editor:** Contains the following SQL code:

```
1 select count(1) from users
```
- Data Output Tab:** Active, showing the result of the query:

	count	bigint
1	1987869	
- Messages and Notifications:** Tabs below Data Output.
- File and Database Management Icons:** At the bottom left.

During the dataset importation process to the Azure database, some issues were faced and solved, including character encoding problems, language errors, database connectivity, and other related

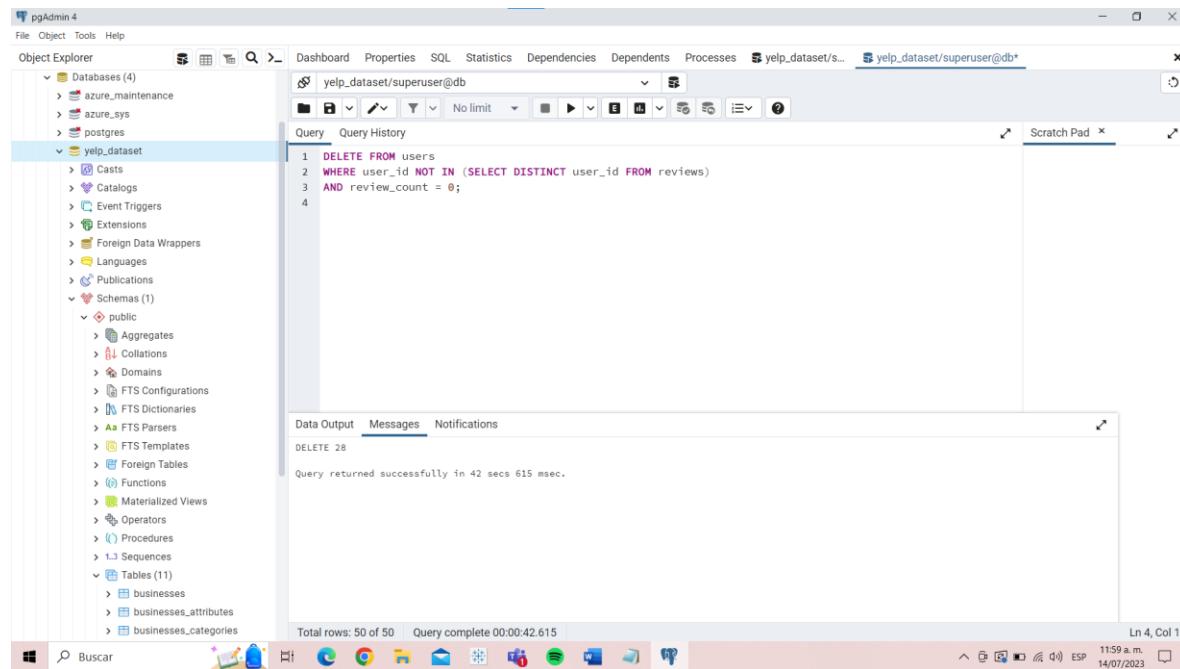
issues. It was a terrific opportunity to learn more about the cloud and how to deal with situations like these.

Dataset cleaning process

Specific queries were run during the cleaning phase to check for duplicate or inconsistent data and remove it from the analysis while maintaining the necessary data to have an accurate analysis and better results. The following queries were executed:

1. Delete inactive users without reviews: Inactive users with no reviews may be considered stale data or abandoned accounts. By removing them, metrics relating to active users become more accurate, and data analysis becomes easier.

```
DELETE FROM users
WHERE user_id NOT IN (SELECT DISTINCT user_id FROM reviews)
AND review_count = 0;
```



The screenshot shows the pgAdmin 4 interface. The left sidebar is the Object Explorer, showing databases (yelp_dataset, azure_maintenance, azure_sys, postgres), casts, catalogs, event triggers, extensions, foreign data wrappers, languages, publications, and a single schema named public containing various objects like aggregates, collations, domains, FTS configurations, FTS dictionaries, FTS parsers, FTS templates, foreign tables, functions, materialized views, operators, procedures, sequences, and tables (businesses, businesses_attributes, businesses_categories). The main area is the Query History tab, which displays the following SQL code:

```
1 DELETE FROM users
2 WHERE user_id NOT IN (SELECT DISTINCT user_id FROM reviews)
3 AND review_count = 0;
4
```

Below the code, the Data Output tab shows the result of the query: "DELETE 28". A message at the bottom states "Query returned successfully in 42 secs 615 msec." and "Total rows: 50 of 50 Query complete 00:00:42.615". The status bar at the bottom right shows "Ln 4, Col 1", "1159 a.m.", "ESP", and the date "14/07/2023".

2. Identify records in attributes, categories, hours, reviews tables that are not associated to

an id presented in the main table “business”: This query is executed to analyse the consistency of the migrated dataset. The mentioned tables should not have records without a “business_id” that is not present in the main table "businesses". If this happens, it must be analysed to determine if it affected the integrity of the migrated data or if the dataset presents inconsistencies since all these features must be related to a business.

OUTPUT:

Businesses_attributes

Properties SQL Dependencies Processes Untitled* yelp_dataset/superuser@dbreviews*

yelp_dataset/superuser@dbreviews

Query History

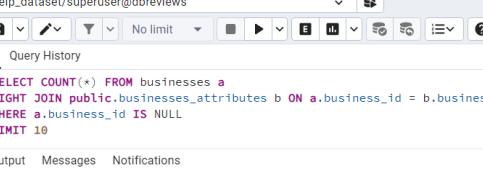
```
1 SELECT COUNT(*) FROM businesses a
2 RIGHT JOIN public.businesses_attributes b ON a.business_id = b.business_id
3 WHERE a.business_id IS NULL
4 LIMIT 10
```

Data Output Messages Notifications

	count	bigint
1	0	

Total rows: 1 of 1 Query complete 00:00:00.481

ES Español (Colombia) Corrección



businesses_categories

Properties SQL Dependencies Processes Untitled* yelp_dataset/superuser@dbreviews*

Query History

```
1 SELECT COUNT(*) FROM businesses a
2 RIGHT JOIN public.businesses_categories b ON a.business_id = b.business_id
3 WHERE a.business_id IS NULL
4 LIMIT 10
```

Data Output Messages Notifications

	count	bigint
1	0	

Total rows: 1 of 1 Query complete 00:00:00.263 ES Español (Colombia) Corrección 40 12:55 AM 2023-07-19

Businesses hours

reviews

Properties SQL Dependencies Processes Untitled* yelp_dataset/superuser@dbreviews*

yelp_dataset/superuser@dbreviews

No limit

Query History

```
1 SELECT COUNT(*) FROM businesses a
2 RIGHT JOIN public.reviews b ON a.business_id = b.business_id
3 WHERE a.business_id IS NULL
4 LIMIT 10
```

Data Output Messages Notifications

	count	bignum
1	0	

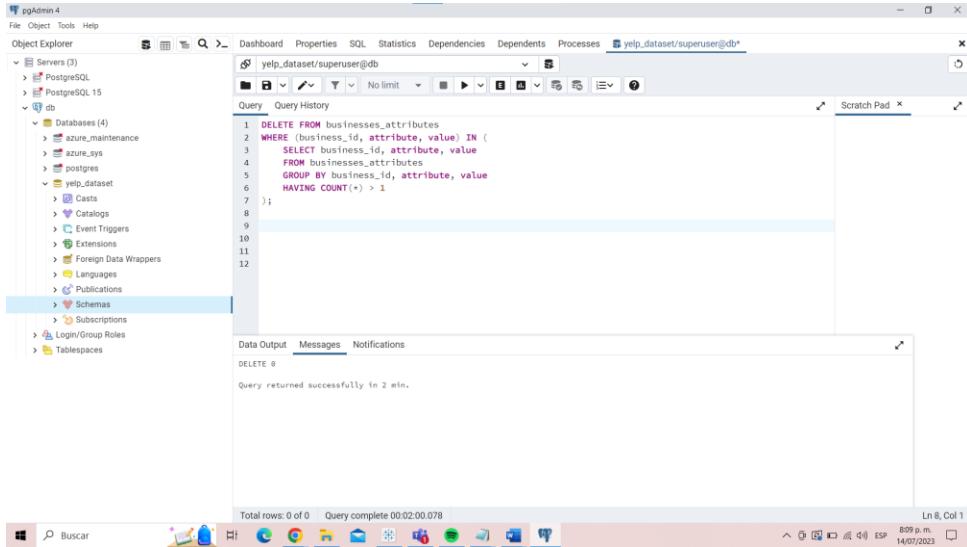
Total rows: 1 of 1 Query complete 00:00:09.579

ES Español (Colombia) Corrección

12:59 AM 2023-07-18

3. Eliminate duplicate business attributes or business categories: If there are duplicate records in the "businesses_attributes" or "businesses_categories" tables with the same combination of business, attribute, categories, and value, it may be redundant data or insert errors. Removing these duplicate attributes helps maintain consistency and prevents confusion when accessing data.

Query business_attributes:



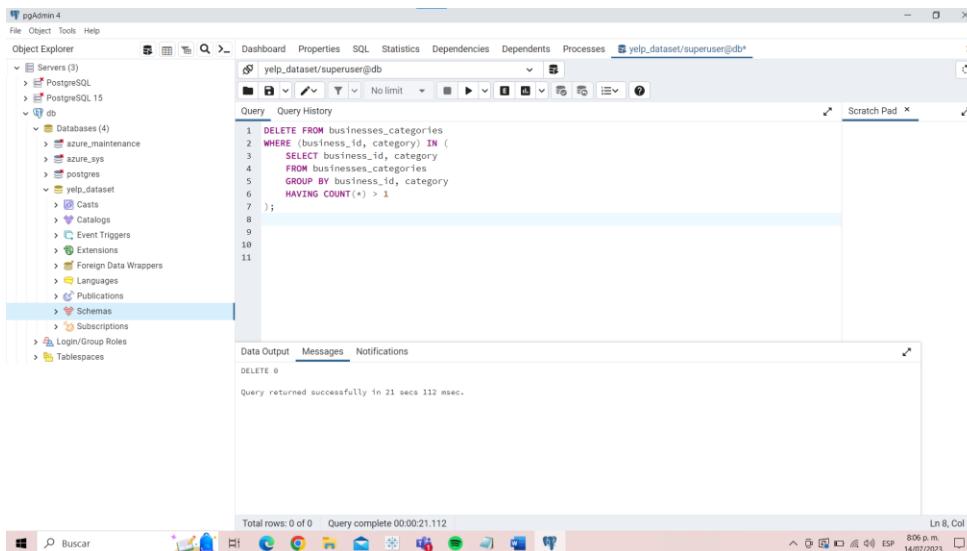
```

pgAdmin 4
File Object Tools Help
Object Explorer Dashboard Properties SQL Statistics Dependencies Dependents Processes yelp_dataset/superuser@db+
Servers (3) PostgreSQL PostgreSQL db
  Databases (4)
    > azure_maintenance
    > azure_sys
    > postgres
    > yelp_dataset
      Casts Catalogs Event Triggers Extensions Foreign Data Wrappers Languages Publications Schemas Subscriptions Login/Group Roles Tablespaces
Query Query History
1 DELETE FROM businesses_attributes
2 WHERE (business_id, attribute, value) IN (
3   SELECT business_id, attribute, value
4     FROM businesses_attributes
5       GROUP BY business_id, attribute, value
6         HAVING COUNT(*) > 1
7
8
9
10
11
12
Data Output Messages Notifications
DELETE 0
Query returned successfully in 2 min.

Total rows: 0 of 0 Query complete 00:02:00.078 Ln 8, Col 1
8:09 p.m. 14/07/2023

```

Query businesses_categories:



```

pgAdmin 4
File Object Tools Help
Object Explorer Dashboard Properties SQL Statistics Dependencies Dependents Processes yelp_dataset/superuser@db+
Servers (3) PostgreSQL PostgreSQL db
  Databases (4)
    > azure_maintenance
    > azure_sys
    > postgres
    > yelp_dataset
      Casts Catalogs Event Triggers Extensions Foreign Data Wrappers Languages Publications Schemas Subscriptions Login/Group Roles Tablespaces
Query Query History
1 DELETE FROM businesses_categories
2 WHERE (business_id, category) IN (
3   SELECT business_id, category
4     FROM businesses_categories
5       GROUP BY business_id, category
6         HAVING COUNT(*) > 1
7
8
9
10
11
Data Output Messages Notifications
DELETE 0
Query returned successfully in 21 secs 112 msec.

Total rows: 0 of 0 Query complete 00:00:21.112 Ln 8, Col 1
8:06 p.m. 14/07/2023

```

After analysing the results of each last query and checking deeply the type and format of data, it was found to be a consistent database with a little "dirt" in just some fields. The information migrated is reliable and convenient for obtaining accurate results.

Dataset analysis and insights

At this stage of the process, a brainstorming session was carried out with the team members to establish specific objectives to obtain significant information when analysing the data. In the first instance, general queries were made that led to the final queries that are shared at this stage of the documentation, with their respective descriptions.

For the data analysis, it was also established that the analysis would focus only on the businesses that are categorised as "Restaurants" in the database.

1. Finding the number of businesses that are cataloged as restaurants: it was executed a query to know the total number of restaurants in the database, since it is the category chosen for this project.

SQL:

```
SELECT COUNT(*)  
FROM businesses_categories  
WHERE category = 'Restaurants';
```

```

SELECT COUNT(*)
  FROM businesses_categories
 WHERE category = 'Restaurants';

```

	count	bigint
1	52268	

Total rows: 1 of 1 Query complete 00:00:00.743

Results: A total of 52,268 businesses categorized as 'Restaurants' were found in the database.

2. Find the 10 cities with the most restaurants and their respective average star rating. With this query, the goal is to find out in which cities YELP has the largest number of registered restaurants and to have an approach to user satisfaction in those cities based on star ratings.

SQL:

```

SELECT b.city||', '||b.state_code AS city,
       COUNT(*) AS restaurant_count,
       ROUND(CAST(AVG(b.stars) AS NUMERIC),2) AS average_stars
  FROM businesses AS b
 JOIN businesses_categories AS bc ON b.business_id = bc.business_id
 WHERE bc.category = 'Restaurants'
 GROUP BY b.city||', '||b.state_code
 ORDER BY restaurant_count DESC
 LIMIT 10;

```

The screenshot shows the pgAdmin 4 interface. On the left is a tree view of database objects under the 'yelp_dataset' schema. The main area contains a SQL query window with the following code:

```

1 SELECT b.city||', '||b.state_code AS city,
2       COUNT(*) AS restaurant_count,
3       ROUND(CAST(AVG(b.stars) AS NUMERIC),2) AS average_stars
4   FROM businesses AS b
5   JOIN businesses_categories AS bc ON b.business_id = bc.business_id
6 WHERE bc.category = 'Restaurants'
7 GROUP BY b.city||', '||b.state_code
8 ORDER BY restaurant_count DESC
9
10 LIMIT 10;

```

Below the query window is a table showing the results of the query:

	city	restaurant_count	average_stars
1	Philadelphia, PA	5852	3.56
2	Tampa, FL	2960	3.60
3	Indianapolis, IN	2862	3.49
4	Nashville, TN	2502	3.56
5	Tucson, AZ	2466	3.46
6	New Orleans, ...	2259	3.74
7	Edmonton, AB	2166	3.50
8	Saint Louis, MO	1790	3.50
9	Reno, NV	1286	3.58
10	Boise, ID	847	3.59

The status bar at the bottom right indicates the query completed at 00:00:00.184 on 2023-07-18 at 11:24 AM.

Results: After this query, it was identified that the city with the largest number of restaurants is Philadelphia with a total of 5,852 businesses in this category, with an approximate average of 3.56 stars. It was also found that 7 of the 10 cities with the largest number of restaurants have more than 2,000 business units of this type.

3. Find the restaurants with the highest stars: This query allows the restaurants with the highest rating (5 stars) to be identified and displayed on the front-end and promoted on the front page.

This helps highlight businesses that offer a great dining experience and attract potential customers looking for high-quality options.

SQL:

```
SELECT b.city||', '||b.state_code AS city,
       b.name, b.stars
  FROM businesses AS b
 JOIN businesses_categories AS bc ON b.business_id = bc.business_id
 WHERE bc.category = 'Restaurants'
   AND b.stars = 5
 ORDER BY city, name ASC;
```

	city text	name character varying	stars double precision
1	Abington, PA	2 Fat Dogs	5
2	Abington, PA	Emoji Sushi & Teriyaki	5
3	Alton, IL	Shake Rattle & Roll Drive In	5
4	Ambler, PA	Backyard Beans Coffee - Ambler	5
5	Ambler, PA	Caffé Maida	5
6	Antioch, TN	Get your Grub on	5
7	Antioch, TN	Mr & Mrs Empanada	5
8	Antioch, TN	Shineworthy Tea	5
9	Apollo Beach, FL	Faith Latin Cuisine	5
10	Apollo Beach, FL	The Mullet Shack	5

Results: After implementing the query, 1,553 restaurants were found with a high rating (5 stars)

by their users.

4. Identify businesses with the worst ratings (<=2) and analyze reviews to identify recurring

issues: By identifying the businesses with the worst ratings, Yelp can identify patterns and specific areas where businesses can improve. This may include aspects such as the quality of the products or services, customer service, cleanliness, waiting time, among others.

This analysis is important because the company is looking at a window of opportunity to expand and diversify the business to increase revenue, and this could help decide if the company could choose to offer a custom assessment to the lowest rated restaurant to identify weaknesses and advise them to improve and achieve better acceptance by users. This advice will be personalized and will be offered at a low cost to any restaurant that wishes to improve its reputation.

Query 1a:

```
SELECT b.name, b.stars, r.review_text
  FROM businesses AS b
  JOIN reviews AS r ON b.business_id = r.business_id
  JOIN businesses_categories AS bc ON b.business_id = bc.business_id
 WHERE bc.category = 'Restaurants'
   AND b.stars <= 2;
```

	name	stars	review_text
1	McDonald's	1	Why does it take sooooo long to get food inside? Or
2	McDonald's	1	Service is slow most times but it's McDonalds so I
3	McDonald's	1	This is the slowest place ever! You can't even cont
4	McDonald's	1	I visited this location Saturday, 6/1/19 about 9:30a
5	McDonald's	1	One thing is certain, McDonalds employees are the
6	McDonald's	1	This place is super slow and rude crew and have s
7	McDonald's	1	Once again sady proof that some people don't de
8	McDonald's	1	***ROACHES*** ***Infestation***
9	McDonald's	1	***ROACHES*** ***Infestation***
10	McDonald's	1	Yo, I ordered ONE thing. Grilled chicken sandwich!
11	McDonald's	1	No indoor seating and no use of restroom. We wer
12	McDonald's	1	I don't know why I even bother stopping. I mean yo
13	McDonald's	1	This restaurant deserves ZERO stars for customer
14	McDonald's	1	Wow. You would think that fast food restaurants w
15	McDonald's	1	I'm a DoorDash driver, I recently got a job to pick u

Query 1b: Identify the number of restaurants with the worst ratings (<=2)

The screenshot shows the pgAdmin 4 interface. The left sidebar displays the schema browser with the 'public' schema selected, showing various objects like Aggregates, Collations, Domains, FTS Configurations, FTS Dictionaries, FTS Parsers, FTS Templates, Foreign Tables, Functions, Materialized Views, Operators, Procedures, and Sequences. Below this, the 'Tables (13)' section is expanded, showing tables such as _aborted, _backup, businesses, reviews, and users. The main area contains a SQL query window with the following code:

```

1 SELECT x.name, COUNT(x.review_text) as review_count
2 FROM (
3     SELECT b.name, b.stars, r.review_text
4         FROM businesses AS b
5             JOIN reviews AS r ON b.business_id = r.business_id
6             JOIN businesses_categories AS bc ON b.business_id = bc.business_id
7     WHERE bc.category = 'Restaurants'
8     AND b.stars <= 2)x

```

The results table shows the top 12 restaurants with the lowest star ratings and their review counts:

	name	review_count
1	McDonald's	7080
2	Taco Bell	2661
3	Wendy's	2182
4	Steak 'n Shake	2110
5	Pizza Hut	1823
6	Burger King	1694
7	Popeyes Louisiana Kitchen	1521
8	Buffalo Wild Wings	1448
9	Domino's Pizza	1426
10	Chipotle Mexican Grill	1421
11	KFC	1156
12	IHOP	1060

Total rows: 765 of 765 | Query complete 00:00:00.499 | ES Español (Colombia) | Corrección | 11

Result: 56,079 restaurant reviews were found with a star rating equal to or less than a score of

2. The top 3 of the restaurants with the worst reviews are MacDonald's, Taco Bell, and Wendy's.

To have more tools when analyzing the results of this query, it was executed three additional queries, filtering three keywords: "customer service", "clean" and "waiting time" grouping for restaurant and city and for those with more than 10 bad reviews.

Query 1c: Number of reviews related to the issues "customer service", "clean" or "waiting time"

```
SELECT COUNT(*)
FROM businesses AS b
JOIN reviews AS r ON b.business_id = r.business_id
JOIN businesses_categories AS bc ON b.business_id = bc.business_id
WHERE b.stars <= 2 AND bc.category = 'Restaurants' AND (r.review_text LIKE
'%customer service%' OR r.review_text LIKE '%clean%' OR r.review_text LIKE
'%wait%');
```

The screenshot shows the pgAdmin 4 interface. The left sidebar displays the Object Explorer with various database objects like public, Aggregates, Collations, Domains, FTS Configurations, FTS Dictionaries, FTS Parsers, FTS Templates, Foreign Tables, Functions, Materialized Views, Operators, Procedures, Sequences, and Tables (13). The central pane shows the SQL editor with the query from above. The Data Output tab shows the result: count bigint, with a value of 19667. The status bar at the bottom indicates the query was completed at 00:00:14.829.

count	bigint
1	19667

Result: 19,667 of 56,079 restaurant reviews that represents 35% were found with the issue "Customer services" or "Clean" or "Wait".

Query 2:

```
SELECT b.city||', '||b.state_code AS city, b.name,
       COUNT(*) reviews,
       ROUND(CAST(AVG(b.stars) AS NUMERIC),2) AS average_rating
  FROM businesses AS b
 JOIN reviews AS r ON b.business_id = r.business_id
 JOIN businesses_categories AS bc ON b.business_id = bc.business_id
 WHERE bc.category = 'Restaurants'
       AND r.review_text LIKE '%customer service%'
 GROUP BY b.city||', '||b.state_code, b.name
 HAVING COUNT(*) > 10 -- with more than 10 reviews
       AND AVG(b.stars) <= 2
```

```
ORDER BY average_rating;
```

The screenshot shows a database interface with a query editor and a results grid. The query is:

```

5   JOIN reviews AS r ON b.business_id = r.business_id
6   JOIN businesses_categories AS bc ON b.business_id = bc.business_id
7   WHERE bc.category = 'Restaurants'
8   AND r.review_text LIKE '%customer service%'
9   GROUP BY b.city||', '||b.state_code, b.name
10  HAVING COUNT(*) > 10 -- with more than 10 reviews
11  AND AVG(b.stars) <= 2
12  ORDER BY average_rating;
13

```

The results grid displays the following data:

	city text	name character varying	reviews bigint	average_rating numeric
1	Tucson, AZ	KFC	19	1.24
2	St. Louis, MO	McDonald's	42	1.37
3	Goleta, CA	Jack in the Box	11	1.50
4	Tampa, FL	El Patron Pizzeria	17	1.50
5	Wesley Chapel, FL	Popeyes Louisiana Kitchen	13	1.50
6	Nashville, TN	McDonald's	25	1.50
7	Sahuarita, AZ	McDonald's	16	1.50
8	Philadelphia, PA	McDonald's	43	1.55
9	Reno, NV	McDonald's	37	1.57
10	Nashville, TN	Wendy's	13	1.58
11	New Orleans, LA	Pizza Hut	19	1.61
12	Tampa, FL	Popeyes Louisiana Kitchen	21	1.64
13	Philadelphia, PA	Wendy's	12	1.67
14	Sparks, NV	Jack in the Box	11	1.68
15	Philadelphia, PA	Pete's Pizza	11	1.68
16	Indianapolis, IN	McDonald's	35	1.70

Total rows: 58 of 58 Query complete 00:01:19.396 ES Español (Colombia) Corrección 1

Results: 58 restaurants with a star rating less than or equal to 2 were found, where "customer service" is mentioned as at least 10 reviews.

Query 3:

```

SELECT b.city||', '||b.state_code AS city, b.name,
       COUNT(*) reviews,
       ROUND(CAST(AVG(b.stars) AS NUMERIC),2) AS average_rating
  FROM businesses AS b
 JOIN reviews AS r ON b.business_id = r.business_id
 JOIN businesses_categories AS bc ON b.business_id = bc.business_id
 WHERE bc.category = 'Restaurants'
   AND r.review_text LIKE '%clean%'
 GROUP BY b.city||', '||b.state_code, b.name
 HAVING COUNT(*) > 10 -- with more than 10 reviews
   AND AVG(b.stars) <= 2

```

ORDER BY average_rating;

	city text	name character varying	reviews bigint	average_rating numeric
1	Tucson, AZ	KFC	12	1.25
2	St. Louis, MO	McDonald's	44	1.36
3	Metairie, LA	McDonald's	11	1.45
4	Upper Darby, PA	IHOP	11	1.50
5	Indianapolis, IN	Clarion Hotel	11	1.50
6	Reno, NV	McDonald's	47	1.53
7	Philadelphia, PA	McDonald's	44	1.56
8	Tampa, FL	Burger King	21	1.76
9	Saint Louis, MO	Chill's	20	1.78
10	Tucson, AZ	Burger King	20	1.83
11	Indianapolis, IN	Buffalo Wild Wings	15	1.87
12	Gretna, LA	McDonald's	15	1.87
13	Sparks, NV	Taco Bell	16	1.88
Total rows: 45 of 45		Query complete 00:01:30.912	ES Español (Colombia)	Corrección

Result: 45 restaurants with a star rating less than or equal to 2 were found, where "clean" is mentioned as at least 10 reviews.

Query 4:

```
SELECT b.city||', '||b.state_code AS city, b.name,
       COUNT(*) reviews,
       ROUND(CAST(AVG(b.stars) AS NUMERIC),2) AS average_rating
  FROM businesses AS b
 JOIN reviews AS r ON b.business_id = r.business_id
 JOIN businesses_categories AS bc ON b.business_id = bc.business_id
 WHERE bc.category = 'Restaurants'
   AND r.review_text LIKE '%clean%'
 GROUP BY b.city||', '||b.state_code, b.name
```

```
HAVING COUNT(*) > 10 -- with more than 10
reviews
```

```
AND AVG(b.stars) <= 2
ORDER BY average_rating;
```

The screenshot shows a database interface with a query editor and a results viewer.

Query Editor:

```

1 SELECT b.city||', '||b.state_code AS city, b.name,
2      COUNT(*) reviews,
3      ROUND(CAST(AVG(b.stars) AS NUMERIC),2) AS average_rating
4  FROM businesses AS b
5  JOIN reviews AS r ON b.business_id = r.business_id
6  JOIN businesses_categories AS bc ON b.business_id = bc.business_id
7 WHERE bc.category = 'Restaurants'
8 AND r.review_text LIKE '%wait%'
9 GROUP BY b.city||', '||b.state_code, b.name
10 HAVING COUNT(*) > 10 -- with more than 10 reviews
11 AND AVG(b.stars) <= 2
12 ORDER BY average_rating;
13

```

Data Output:

	city text	name character varying	reviews bigint	average_rating numeric
1	Chalmette, LA	Burger King	14	1.00
2	Belleville, IL	Steak 'n Shake	18	1.00
3	Sun Valley, NV	McDonald's	14	1.00
4	Camden, NJ	McDonald's	11	1.23
5	Tucson, AZ	KFC	70	1.31
6	Upper Darby, PA	McDonald's	14	1.39
7	St. Louis, MO	McDonald's	181	1.40
8	Metairie, LA	McDonald's	33	1.44
9	Metairie, LA	Burger King	19	1.47
10	Greenfield, IN	McDonald's	14	1.50
11	Goleta, CA	Jack in the Box	15	1.50
12	Wesley Chapel, FL	Popeyes Louisiana Kitchen	36	1.50
13	Philadelphia, PA	McDonald's	117	1.50

Total rows: 296 of 296 Query complete 00:01:31.284 ES Español (Colombia) Corrección 1

Icons at the bottom include: Google, Microsoft Edge, Microsoft Word, Microsoft Excel, Microsoft Powerpoint, Microsoft OneDrive, and a blue circular icon with a white 'S'.

Result: 296 restaurants with a star rating less than or equal to 2 were found, where "wait" is mentioned as at least 10 reviews.

5. Find the most influential users (largest number of followers) and collaborate with them

to promote businesses: By identifying the most influential users, Yelp can establish strategic collaborations with them to promote their services and highlight certain businesses or events. These

influential users can have a significant impact on the opinion and decisions of their followers, which could increase Yelp's visibility and reputation.

SQL:

```
SELECT u.name, u.fans
FROM users AS u
ORDER BY u.fans DESC
LIMIT 10;
```

The screenshot shows the pgAdmin 4 interface. On the left, the Object Explorer tree shows databases like 'yelp_dataset' selected. The main pane displays a query window with the following SQL code:

```
1 SELECT u.name, u.fans
2 FROM users AS u
3 ORDER BY u.fans DESC
4 LIMIT 10;
```

Below the query window is a table titled 'Data Output' showing the results:

	name	fans
1	Mike	12497
2	Katie	3642
3	Fox	3493
4	Richard	3243
5	Daniel	3138
6	Jessica	2627
7	Ruggy	2547
8	Megan	2451
9	Emi	2424
10	Peter	2388

At the bottom of the pgAdmin window, it says 'Total rows: 10 of 10 Query complete 00:00:40.608'. The system tray at the bottom right shows the date and time as '10/06 a.m. 15/07/2023'.

Result: The 10 users with the largest number of fans were found, among these the one with the most followers has a total of 12497, approximately 29% more than the user with the second most fans.

All these users have more than 2,000 followers, a significant number of fans.

6. Find the most active users and reward them with special promotions to build their loyalty:

Knowing the most active users, Yelp can create customer loyalty strategies for them, such as rewarding these users with special promotions, thus strengthening the loyalty of those users towards the platform. By receiving exclusive incentives, these users can feel valued and recognized, which increases their satisfaction and motivates them to continue using Yelp as their preferred platform for finding and reviewing businesses.

SQL:

```
SELECT u.name, u.review_count
FROM users AS u
ORDER BY u.review_count DESC
LIMIT 100;
```

The screenshot shows the pgAdmin 4 interface. The left sidebar displays the Object Explorer with servers, databases, and tables. The main area shows a query editor with the following SQL code:

```
1 SELECT u.name, u.review_count
2 FROM users AS u
3 ORDER BY u.review_count DESC
4 LIMIT 100;
```

Below the query, the Data Output tab shows the results:

	name	review_count
1	Fox	17473
2	Victor	16978
3	Bruce	16567
4	Shila	12868
5	Kim	9941
6	Nijole	8363
7	Vincent	8354
8	George	7738
9	Kenneth	6766
10	Jennifer	6679
11	Sunil	6459
12	Eric	5887
13	Ed	5800
14	Rob	5511

Total rows: 100 of 100 Query complete 00:00:30.674

Return: This was the result obtained for the 100 users with the largest number of reviews, the first of which is called Fox and has 17,473 reviews.

After implementing the improvements in its services, YELP will be able to recommend the most active users visit these restaurants and rate them in order to improve their score and start a loyalty programme where these same people attract more people.

Spark

Running some jobs using PySpark

This section will provide a summarized report of the PySpark jobs that were run for this project,

The purpose of this section is to document the results of running PySpark jobs on our set of CSV files.

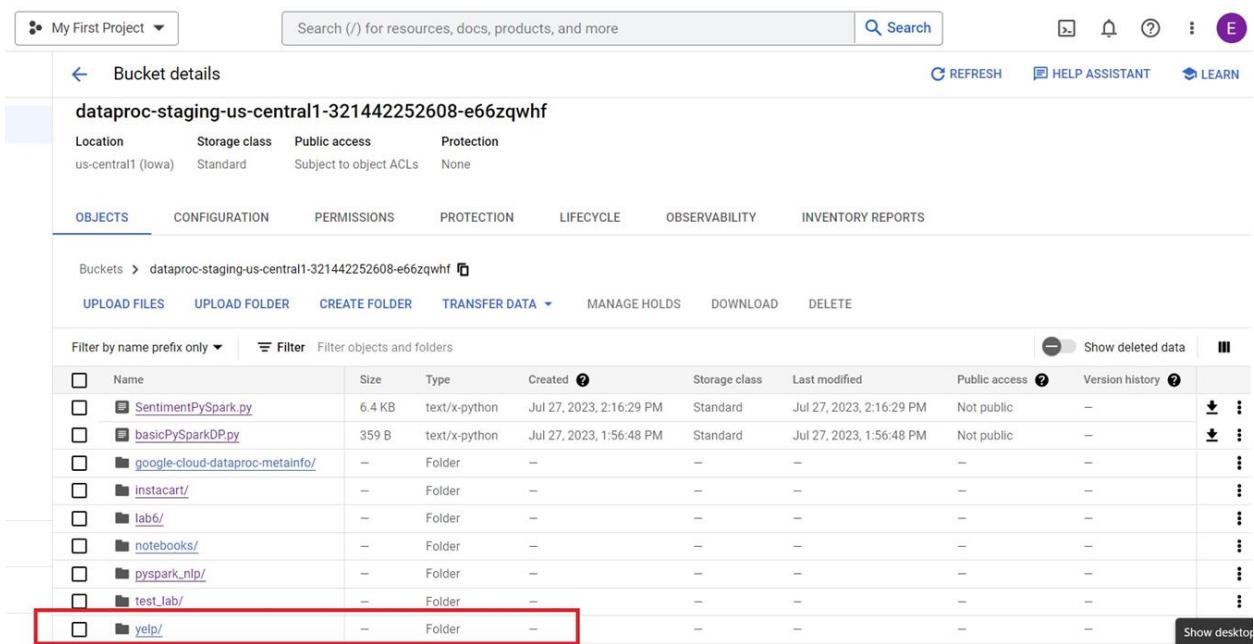
The data in the CSV files consists of the cleaned dataset as of the previous step's output, and the

PySpark jobs' goal was to analyze the data and identify trends.

Methodology

The following steps were taken to run the PySpark jobs:

Firstly A specific working folder was created.



The screenshot shows the Google Cloud Storage 'Bucket details' interface. At the top, there is a navigation bar with 'My First Project' dropdown, search bar, and various icons. Below the navigation bar, the bucket name 'dataproc-staging-us-central1-321442252608-e66zqwhf' is displayed. Under the 'OBJECTS' tab, a list of objects is shown, including several folders like 'SentimentPySpark.py', 'basicPySparkDP.py', and 'yelp/'. The 'yelp/' folder is highlighted with a red border. The table columns include Name, Size, Type, Created, Storage class, Last modified, Public access, and Version history. At the bottom right of the table, there is a 'Show desktop' button.

Name	Size	Type	Created	Storage class	Last modified	Public access	Version history
SentimentPySpark.py	6.4 KB	text/x-python	Jul 27, 2023, 2:16:29 PM	Standard	Jul 27, 2023, 2:16:29 PM	Not public	-
basicPySparkDP.py	359 B	text/x-python	Jul 27, 2023, 1:56:48 PM	Standard	Jul 27, 2023, 1:56:48 PM	Not public	-
google-cloud-dataproc-metainfo/_	-	Folder	-	-	-	-	-
instacart/_	-	Folder	-	-	-	-	-
lab6/_	-	Folder	-	-	-	-	-
notebooks/_	-	Folder	-	-	-	-	-
pyspark_nlp/_	-	Folder	-	-	-	-	-
test_lab/_	-	Folder	-	-	-	-	-
yelp/_	-	Folder	-	-	-	-	-

For the second step, The CSV files were uploaded to the working folder.

Bucket details for dataproc-staging-us-central1-321442252608-e66zqwhf

Name	Size	Type	Created	Storage class
categories.csv	23.4 MB	text/csv	Aug 11, 2023, 4:37:49 PM	Standard
rest_categories.csv	2.4 KB	text/csv	Aug 11, 2023, 9:49:27 PM	Standard
reviews.csv	4.6 GB	text/csv	Aug 11, 2023, 6:08:31 PM	Standard
users.csv	148.3 MB	text/csv	Aug 11, 2023, 4:38:36 PM	Standard

And next, A PySpark job was created to process the CSV files.

Submit a job

Job ID *	job-adfe81dd
Region *	us-central1
Cluster *	cluster-316a
Job type *	PySpark
Main python file *	gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/reviews_stats.py
Additional python files	
Jar files	
Files	

My First Project ▾

Search (/) for resources, docs, products, and more

Search

7 ? E

Submit a job

Arguments

gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/data/

gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/results

Additional arguments to pass to the main class. Press Return after each argument.

Max restarts per hour

Leave blank if you don't want to allow automatic restarts on job failure. [Learn more](#)

Properties [?](#)

+ ADD PROPERTY

Labels

+ ADD LABEL

SUBMIT CANCEL

EQUIVALENT REST

After that, The PySpark job was launched.

My First Project ▾

Search (/) for resources, docs, products, and more

Search

7 ? E

Job details

CLONE DELETE STOP REFRESH

Job ID: job-adfe81dd

Job UUID: caffe38c-029f-4645-93fe-e4cd4f53942d

Type: Dataproc Job

Status: Starting

MONITORING CONFIGURATION

The charts below represent the metrics from the cluster this job ran on, scoped to the time that this job was running. It is possible for more than one job to run on a cluster at a time, so these metrics may not reflect this job's resource usage accurately. Metrics for a job may lag behind the job run by several minutes.

Output LINE WRAP: OFF

Press Alt+F1 for Accessibility Options.

The PySpark job was monitored until it was completed successfully.

Job is running

The screenshot shows the 'Job details' page for a Dataflow job named 'job-adfe81dd'. The status is 'Running'. The output section displays log entries:

```

Press Alt+F1 for Accessibility Options.
23/08/14 20:52:37 INFO ResourceUtils: Unable to find 'resource-types.xml'.
23/08/14 20:52:37 INFO YarnClientImpl: Submitted application application_1691964618340_0039
23/08/14 20:52:39 INFO DefaultIoRMFailoverProxyProvider: Connecting to ResourceManager at cluster-316a-m.us-central1-c.c.studied-indexer-387116.internal./10.128.0.4:8030
23/08/14 20:52:41 WARN GfsStorageStatistics: Detected potential high latency for operation op_get_file_status. latencyMs=235; previousMaxLatencyMs=0; operationCount=1; context=gs://dataproc-temp-1691964618340-0039
23/08/14 20:52:41 INFO GoogleCloudStorageImpl: Ignoring exception of type GoogleJsonResponseException; verified object already exists with desired state.
23/08/14 20:52:41 WARN GfsStorageStatistics: Detected potential high latency for operation op_mkdirs. latencyMs=198; previousMaxLatencyMs=0; operationCount=1; context=gs://dataproc-temp-1691964618340-0039
Defining spark sql context...OK
  
```

Job successfully completed

The screenshot shows the 'Job details' page for the same Dataflow job, now with status 'Succeeded'. The output section displays a large amount of log entries, starting with:

```

| Ambler| Fast Food|3.805825242718467|3.0228458309421407|0.5132973519093951|
| Ambler| Sports Bars|2.9761904761904763|3.4531277614603235|0.5132973519093951|
| Antioch| Fast Food|2.0138504155124655|3.0228458309421407|0.5132973519093951|
| Antioch| Sports Bars|2.9881656804733727|3.4531277614603235|0.5132973519093951|
| Apollo Beach| Fast Food| 3.894736842105263|3.0228458309421407|0.5132973519093951|
| Apollo Beach| Sports Bars| 3.5466666666666667|3.4531277614603235|0.5132973519093951|
| Arabi| Fast Food| 4.705882352941177|3.0228458309421407|0.5132973519093951|
| | Fast Food|1.8611111111111112|3.0228458309421407|0.5132973519093951|
| Ardmore| Sports Bars| 3.611940298597463|3.4531277614603235|0.5132973519093951|
| Arnold| Fast Food|2.8973105134474326|3.0228458309421407|0.5132973519093951|
| Arnold| Sports Bars| 3.513157894736842|3.4531277614603235|0.5132973519093951|
| Ashland City| Fast Food| 1.5625|3.0228458309421407|0.5132973519093951|
| Aston| Fast Food|2.4313725490196076|3.0228458309421407|0.5132973519093951|
+-----+
only showing top 20 rows
Getting final results...OK
  
```

The screenshot shows the Data Pipeline interface with a search bar at the top. Below it, a table lists a single job entry:

Job ID	Status	Region	Type	Cluster	Start time	Elapsed time
job-adfe81dd	Succeeded	us-central1	PySpark	cluster-316a	Aug 14, 2023, 4:52:27 PM	13 min 6 sec

A red box highlights the first row of the table.

The results of the PySpark job were checked out.

The screenshot shows the Google Cloud Storage bucket details page for 'dataproc-staging-us-central1-321442252608-e66zqwhf'. It displays the following information:

Location	Storage class	Public access	Protection
us-central1 (Iowa)	Standard	Subject to object ACLs	None

Below this, the 'OBJECTS' tab is selected, showing a list of files and folders:

Name	Size	Type	Created	Storage class	Last modified	Public access	Version history
data/	—	Folder	—	—	—	—	—
results/	—	Folder	—	—	—	—	—
reviews_stats.py	6.1 KB	text/x-python	Aug 14, 2023, 4:48:18 PM	Standard	Aug 14, 2023, 4:48:18 PM	Not public	—

A red box highlights the 'data/' and 'results/' folder entries in the list.

The screenshot shows the 'Bucket details' page for the bucket 'dataproc-staging-us-central1-321442252608-e66zqwhf'. The 'OBJECTS' tab is selected, displaying a list of objects. A red box highlights the first five objects in the list:

Name	Size	Type	Created	Storage class	Last modified	Public access	Version
<u>SUCCESS</u>	0 B	application/octet-stream	Aug 14, 2023, 5:05:29 PM	Standard	Aug 14, 2023, 5:05:29 PM	Not public	-
part-00000-5257d2aa-ec75-4806-9274-8eee9dff1d56-c000.csv	299 B	application/octet-stream	Aug 14, 2023, 5:04:16 PM	Standard	Aug 14, 2023, 5:04:16 PM	Not public	-
part-00000-87119f23-5425-4633-ba44-e0c26156f3d3-c000.csv	9.6 KB	application/octet-stream	Aug 14, 2023, 5:02:59 PM	Standard	Aug 14, 2023, 5:02:59 PM	Not public	-
part-00000-936574ae-6e0c-4e73-80ba-813bc97c8824-c000.csv	43 B	application/octet-stream	Aug 14, 2023, 5:02:03 PM	Standard	Aug 14, 2023, 5:02:03 PM	Not public	-
part-00000-c2145ac6-c6a5-4ed6-acb6-3be86e8c2a10-c000.csv	5.9 KB	application/octet-stream	Aug 14, 2023, 5:03:30 PM	Standard	Aug 14, 2023, 5:03:30 PM	Not public	-
part-00000-f51d97fb-b3c6-4f1f-9... (truncated)	66.8 KB	application/octet-stream	Aug 14, 2023, 5:05:29 PM	Standard	Aug 14, 2023, 5:05:29 PM	Not public	-

The PySpark jobs were successfully done, and the outputs will present the results of the PySpark jobs for use in the next Step: Visualization

The PySpark jobs that were run produced five output files:

- `yelp_results_part-00000-936574ae-6e0c-4e73-80ba-813bc97c8824-c000.csv`: (`qty_categories.csv`)

The output file contains the number of categories that are above the general average rating, whereas the number of categories that are below the general average rating.

- `yelp_results_part-00000-87f19f23-5425-4633-ba44-e0c26156f3d3-c000.csv`: (`avg_categories.csv`)

The output file contains the average rating by category, with the general average rating in the last column.

- `yelp_results_part-00000-5257d2aa-ec75-4806-9274-8eee9dff1d56-c000.csv`:

The output file contains the ten best and worst-rated categories, with the best-rated categories listed first.

- `yelp_results_part-00000-c2145ac6-c6a5-4ed6-acb6-3be86e8c2a10-c000.csv`: (`standard_deviation.csv`)

The output file contains the standard deviation of the ratings by category.

- `yelp_results_part-00000-f51d97fb-b3c6-4f1f-9409-16b7c9a809cc-c000.csv`:

This file contains the average rating by city and category, with the correlation coefficient between the average rating and the general average rating in the last column.

The PySpark jobs provided a lot of interesting insights into the Yelp dataset. The findings from the jobs can be used to inform future marketing and business decisions for businesses on Yelp. For example, businesses in the "Restaurants" category may want to focus on improving their customer service and food quality in order to maintain their high average rating.

The related code for running PySpark Jobs and the outputs file are uploaded to the project [git repository](#). Please check it for any further information.

Findings from the result it would be more understandable by looking at graphs and charts. Let's do the Visualizations Part

Visualization

Technologies:

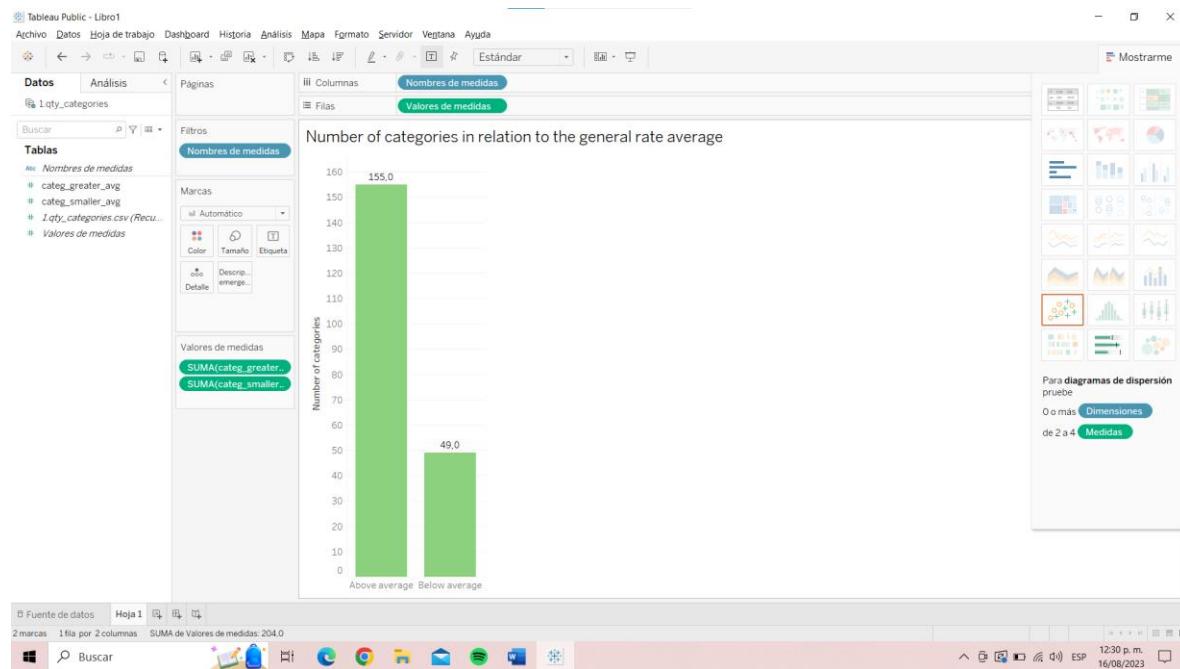
We chose Tableau for our project due to its ability to create interactive visualizations, effective presentations, and data-driven analysis through a wide range of visualizations and data connections, thereby enhancing our informed decision-making.

Also using Matplotlib due the power of ease and quick plot generation

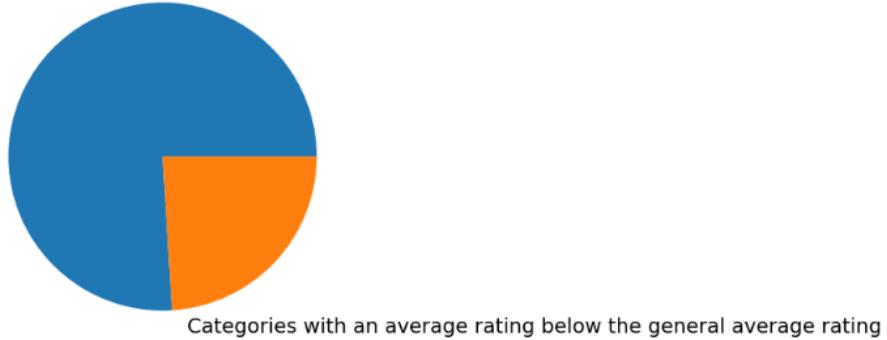
The following are some of the key findings from the PySpark jobs based on the extracted graphs:

Number of categories in relation to the general rate average.

The two graphs reveal that the distribution of restaurant categories is uneven in relation to their average ratings. A significant majority of categories (155) are rated above average, while a much smaller number (49) are rated below average. This indicates that there is a wide range of restaurant quality, with many restaurants performing well above average and a few performing below average.

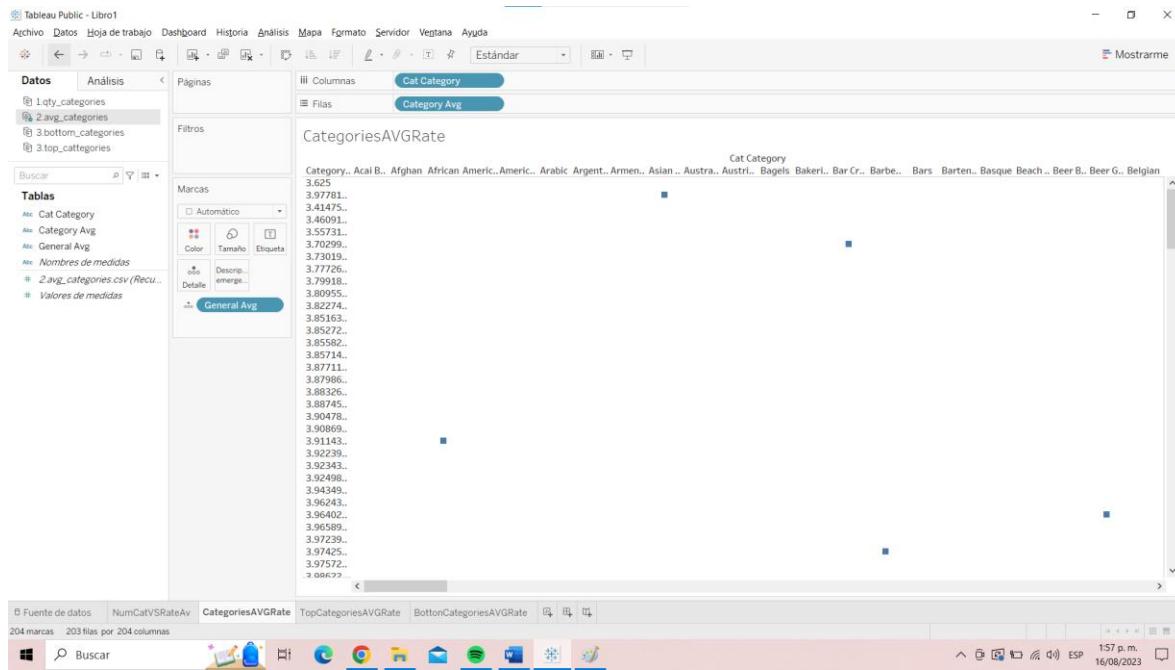


Categories with an average rating above the general average rating

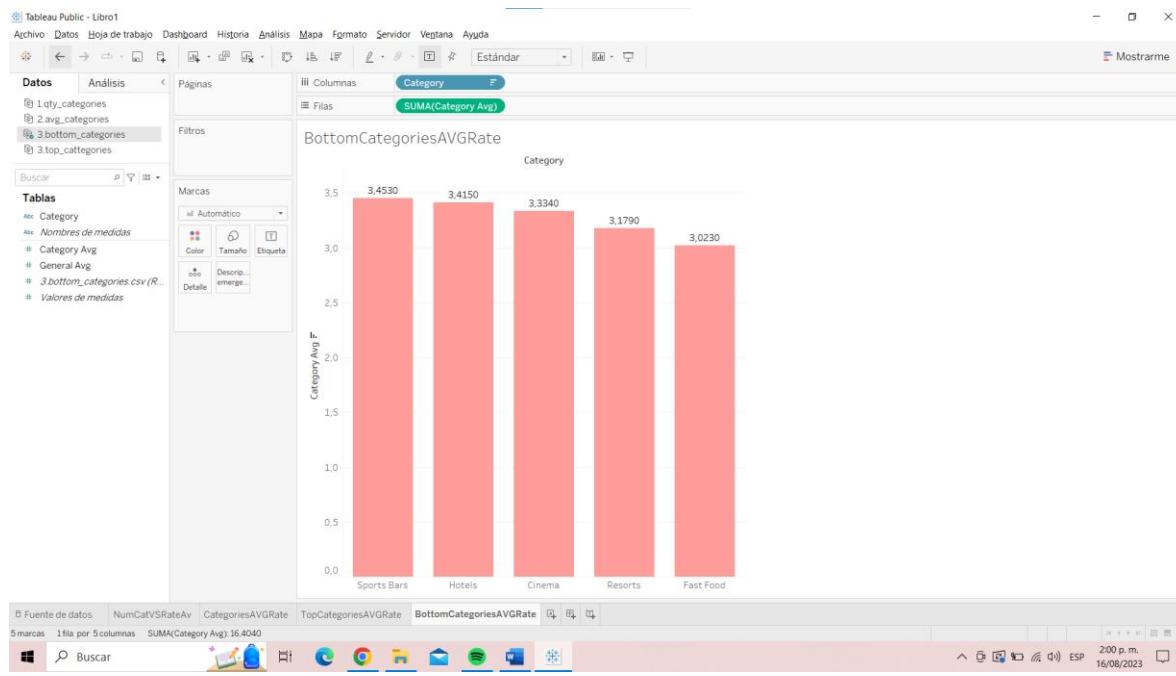
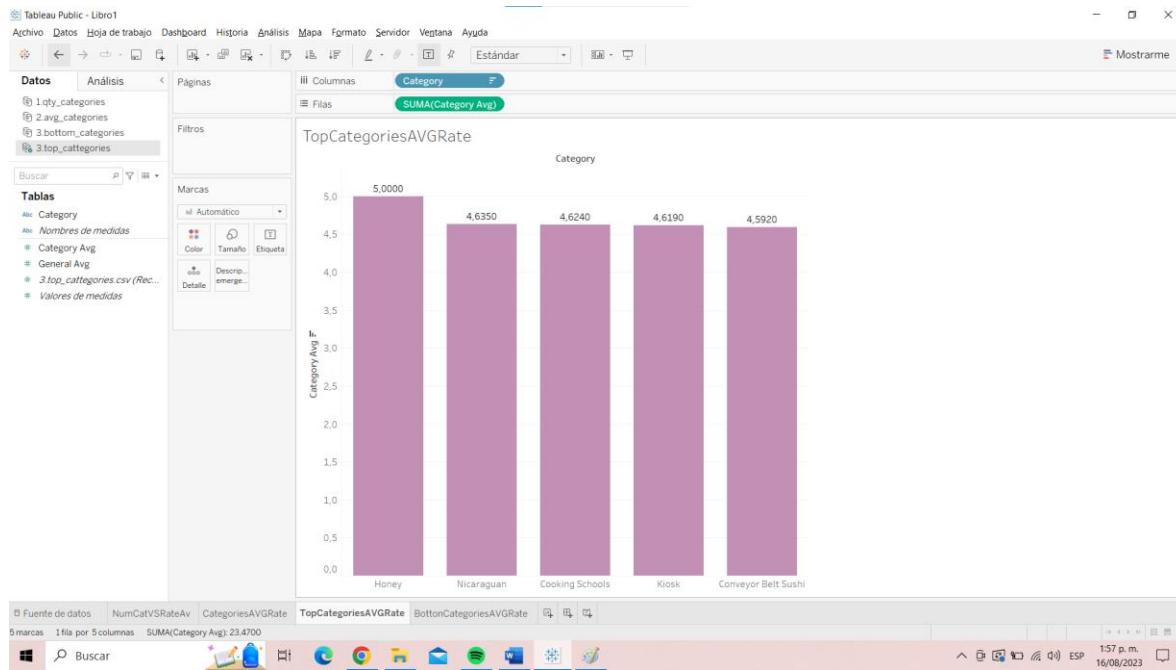


Average Rates of Each Category

At this point of view we wanted to find the average rates for each of the restaurant categories and compare it to the general average.



However, as can be seen in the visualization, since there is so much data (204) the statistics cannot be clearly observed, for this reason we decided to divide the table into two new tables, to be able to graph the top 5 (above the general average of rates) and the bottom 5 (below the general rate average) of the rate averages by category and to be able to observe them more clearly, these are the results:



Average Rates of Each Category

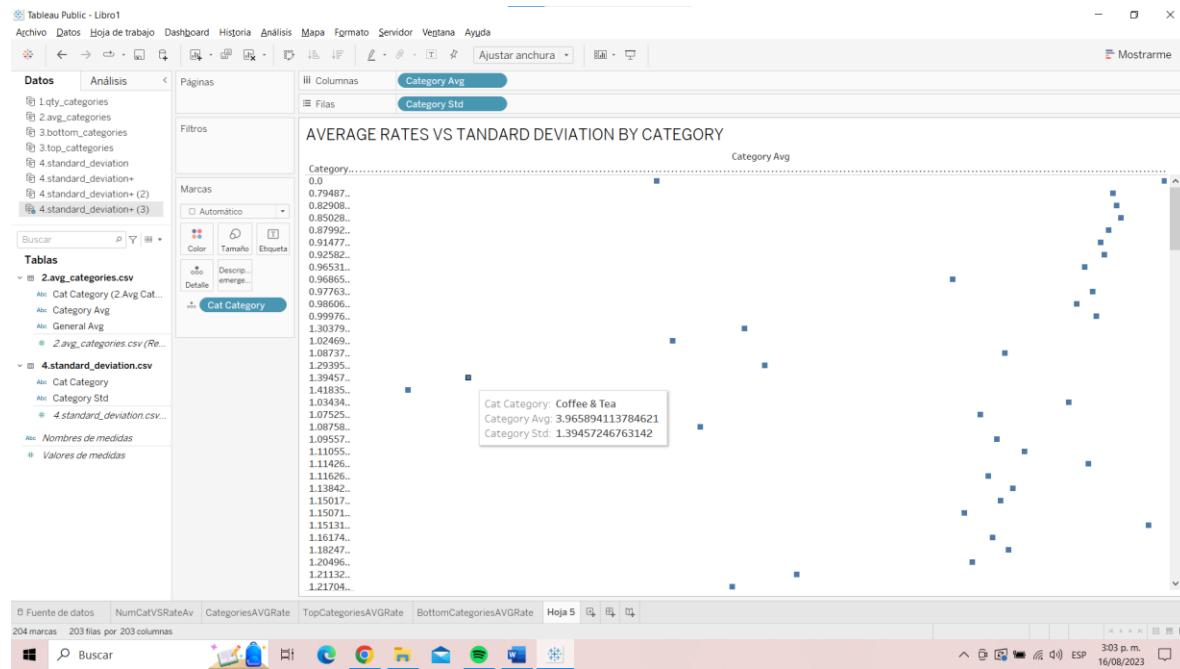
An apparent trend was found that the categories with the highest average have less deviation, and the categories with the lowest average have greater deviation; however, the presence of enough data that do not follow this pattern indicates that the relationship is not necessarily absolute and that there are other factors to consider, for example:

Diversity of Experiences: Categories that are not following the trend may be experiencing a wide range of experiences and opinions. This could indicate that customer experiences within these categories are diverse and that there are factors that drive widely varying opinions.

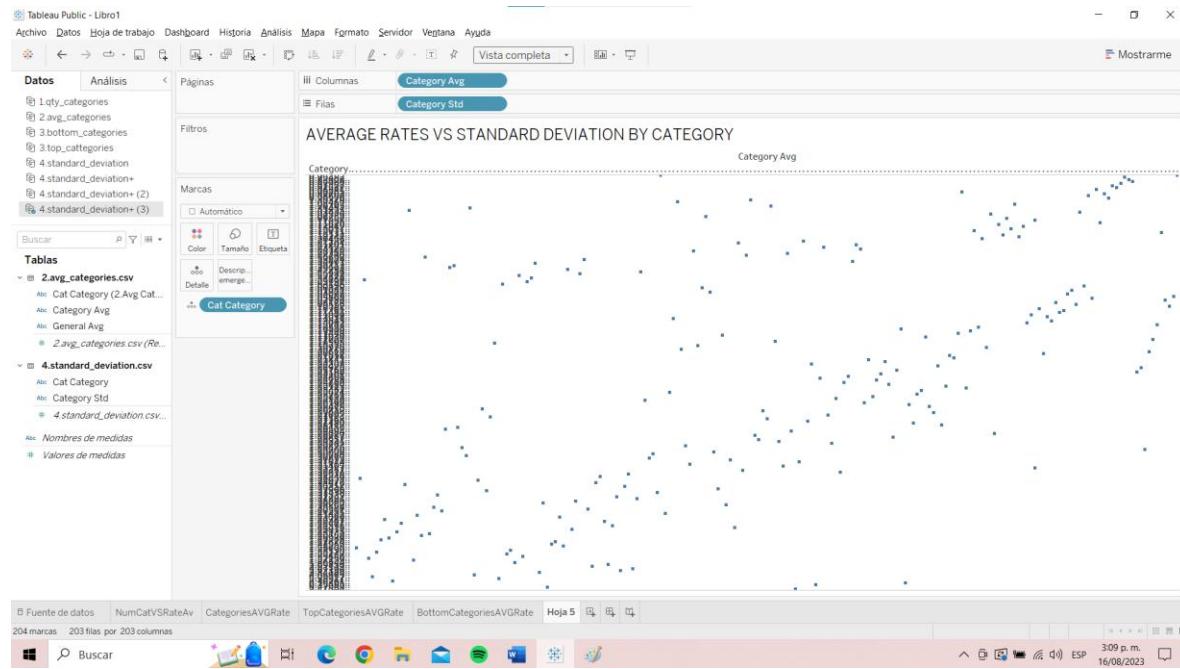
Influencing Factors: Non-trending data may be influenced by unique factors that do not apply to trending categories. There may be particular characteristics of those categories that affect the relationship between the mean and the deviation.

In conclusion, the relationship between the average rate and the standard deviation is complex and cannot be uniformly generalized to all categories. It is essential to take into account the specific context of each category to understand why some follow the trend and others do not.

Average Rates VS. Standard Deviation by Category



A little Zoom Out :



Map Reduce

MapReduce is a programming model and an associated implementation for processing and generating large data sets. It is designed to scale up from single computers to thousands of machines. MapReduce programs are written in Java, Python, or other languages that can be compiled to Java bytecode.

The MapReduce programming model consists of two parts:

- The map function: This function takes a key/value pair as input and produces a set of intermediate key/value pairs as output.
- The reduce function: This function takes a set of intermediate key/value pairs as input and produces a final set of key/value pairs as output.

The benefits of MapReduce include:

Scalability: MapReduce can be scaled up to process very large data sets.

Fault tolerance: MapReduce is fault-tolerant, meaning that it can continue to process data even

Here are some of the use cases of MapReduce:

Analyzing large data sets, such as financial transactions, customer logs, or social media data.

Distributing tasks across a cluster of computers.

Processing data in real time.

Generating reports and summaries.

Storing and retrieving data.

Due this ability of Map Reduce we are going to run some jobs using that to be more familiar with the capabilities

Running some jobs using Map Reduce

A Map Reduce job ran to count the words on top-rated reviews and low created reviews

Counting the words in top-rated reviews and low-rated reviews can be helpful for data analysts in a number of ways. Here are a few examples:

Identifying the most common words in each category. This can help analysts to understand the key themes and topics that are discussed in the reviews. For example, if the most common words in top-rated reviews are "food," "service," and "ambiance," then analysts can infer that these are the factors that are most important to customers when rating a restaurant.

Comparing the word counts between the two categories. This can help analysts to identify any differences in the language used to describe top-rated and low-rated reviews. For example, if the most common words in low-rated reviews are "disappointing" and "unfriendly," then analysts can infer that these are the factors that are most likely to lead to a low rating.

Identifying outliers. If there are any words that are much more common in one category than the other, then these could be considered outliers. Outliers can be helpful for identifying unusual patterns or trends in the data.

Generating reports and visualizations. The word counts can be used to generate reports and visualizations that can be used to communicate the findings of the analysis to stakeholders.

Uploading a new python program to getting reviews from categories which are over and under the general avg rate. This will be the input for map reduce process.

The screenshot shows two pages from the Google Cloud Storage interface:

Bucket details:

- Bucket Name:** dataproc-staging-us-central1-321442252608-e66zqwhf
- Location:** us-central1 (Iowa)
- Storage class:** Standard
- Public access:** Subject to object ACLs
- Protection:** None

OBJECTS (List of objects in bucket):

Name	Type	Size	Created	Storage class	Last modified	Public access	Version history
data/	Folder	—	—	—	—	—	—
results/	Folder	—	—	—	—	—	—
reviews_mapr.py	text/x-python	3.2 KB	Aug 16, 2023, 4:26:30 PM	Standard	Aug 16, 2023, 4:26:30 PM	Not public	—
reviews_stats.py	text/x-python	6.1 KB	Aug 14, 2023, 4:48:18 PM	Standard	Aug 14, 2023, 4:48:18 PM	Not public	—

Object details for reviews_mapr.py:

- LIVE OBJECT**
- VERSION HISTORY**
- Actions:** DOWNLOAD, EDIT METADATA, EDIT ACCESS, DELETE

Overview:

Type	text/x-python
Size	3.2 KB
Created	Aug 16, 2023, 4:26:30 PM
Last modified	Aug 16, 2023, 4:26:30 PM
Storage class	Standard
Custom time	—
Public URL	Not applicable
Authenticated URL	https://storage.cloud.google.com/dataproc-staging-us-central1-321442252608-e66zqwhf/veln/reviews_mapr.py
gsutil URI	gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/reviews_mapr.py

Permissions:

Public access	Not public
---------------	------------

Protection:

Version history	—
-----------------	---

Submitting a new job

Job ID *
job-a877b6a0

Region *
us-central1

Specifies the Cloud Dataproc regional service, which determines what clusters are available.

Cluster *
cluster-316a

Job type *
PySpark

Main python file *
gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/reviews_mapr.py

Can be a GCS file with the gs:// prefix, an HDFS file on the cluster with the hdfs:// prefix, or a local file on the cluster with the file:// prefix

Additional python files

Jar files

Jar files are included in the CLASSPATH. Can be a GCS file with the gs:// prefix, an HDFS file on the cluster with the hdfs:// prefix, or a local file on the cluster with the file:// prefix.

Archive files

Archive files are extracted in the Spark working directory. Can be a GCS file with the gs:// prefix, an HDFS file on the cluster with the hdfs:// prefix, or a local file on the cluster with the file:// prefix. Supported file types: jar, tar, tar.gz, tgz, zip.

Arguments

gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/data/ ↗
gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/results_mapr ↗

Press <Return> to add more arguments

Additional arguments to pass to the main class. Press Return after each argument.

Max restarts per hour

Leave blank if you don't want to allow automatic restarts on job failure. [Learn more ↗](#)

Properties ?

+ ADD PROPERTY

Labels

+ ADD LABEL

SUBMIT CANCEL

Running the job

The screenshot shows the 'Job details' page for a job named 'job-a877b6a0'. The job is a 'Dataproc Job' with a status of 'Running'. The 'Status' field is highlighted with a red box. Below the job details, there are tabs for 'MONITORING' and 'CONFIGURATION'. A note states: 'The charts below represent the metrics from the cluster this job ran on, scoped to the time that this job was running. It is possible for more than one job to run on a cluster at a time, so these metrics may not reflect this job's resource usage accurately. Metrics for a job may lag behind the job run by several minutes.' In the 'Output' section, the logs show:

```

Press Alt+F1 for Accessibility Options.
|     general_avg|
+-----+
|3.9094395932252817|
+-----+
Getting general avg of ratings...OK
Getting avg rates by category...OK

```

Job successfully completed

The screenshot shows the 'Job details' page for the same job, now with a status of 'Succeeded'. The 'Status' field is highlighted with a red box. The 'Output' section displays the full log output:

```

23/08/16 20:34:43 INFO GoogleCloudStorageFileSystem: Successfully repaired 'gs://dataproc-staging-us-central1-321442252608-e66zqwhf/yelp/results_mapr/' directory.
23/08/16 20:34:43 WARN GhfsStorageStatistics: Detected potential high latency for operation op_delete. latencyMs=196; previousMaxLatencyMs=0; operationCount=1; context=gs://dataproc-stagin
23/08/16 20:34:43 WARN GhfsStorageStatistics: Detected potential high latency for operation op_create. latencyMs=107; previousMaxLatencyMs=94; operationCount=2; context=gs://dataproc-stagi
23/08/16 20:34:43 WARN GhfsStorageStatistics: Detected potential high latency for operation stream_write_close_operations. latencyMs=196; previousMaxLatencyMs=0; operationCount=1; context=
23/08/16 20:36:36 WARN GhfsStorageStatistics: Detected potential high latency for operation stream_write_close_operations. latencyMs=440; previousMaxLatencyMs=196; operationCount=2; context
Saving results...OK

Output is complete

```

Results folder

Bucket details for **dataproc-staging-us-central1-321442252608-e66zqwhf**

Location	Storage class	Public access	Protection
us-central1 (Iowa)	Standard	Subject to object ACLs	None

OBJECTS

Buckets > dataproc-staging-us-central1-321442252608-e66zqwhf > **yelp**

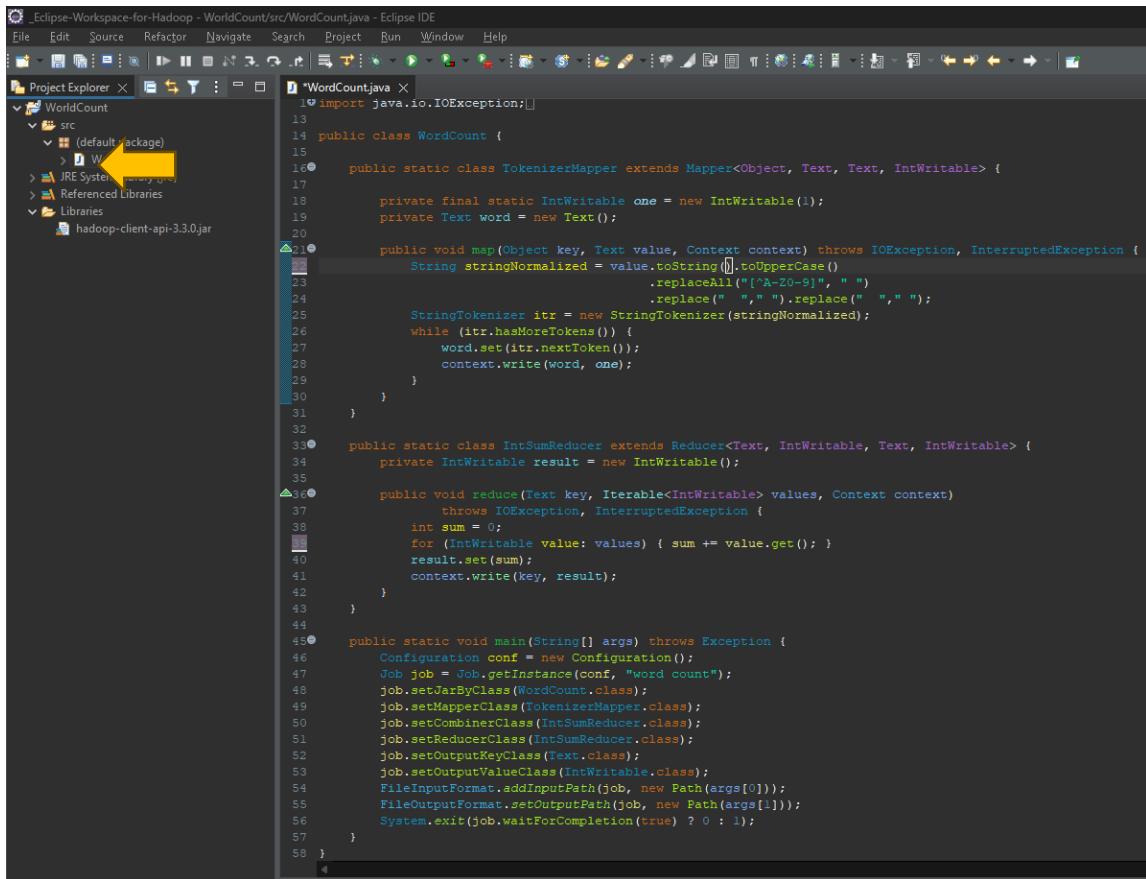
Name	Type	Created	Storage class	Last modified	Public access	Version history
data/	Folder	—	—	—	—	—
results/	Folder	—	—	—	—	—
results_mpr/	Folder	—	—	—	—	—
reviews_mpr.py	text/x-python	Aug 16, 2023, 4:26:30 PM	Standard	Aug 16, 2023, 4:26:30 PM	Not public	Download
reviews_stats.py	text/x-python	Aug 14, 2023, 4:48:18 PM	Standard	Aug 14, 2023, 4:48:18 PM	Not public	Download

Resulting csv files

Bucket details for **dataproc-staging-us-central1-321442252608-e66zqwhf** > **yelp** > **results_mpr**

Name	Type	Created	Storage class	Last modified	Public access	Version history
_SUCCESS	application/octet-stream	Aug 16, 2023, 4:36:36 PM	Standard	Aug 16, 2023, 4:36:36 PM	Not public	Download
part-00000-3079ec3f-c161-4cdd-...	application/octet-stream	Aug 16, 2023, 4:36:14 PM	Standard	Aug 16, 2023, 4:36:14 PM	Not public	Download
part-00000-4ef017bc-917d-42f7-b...	application/octet-stream	Aug 16, 2023, 4:34:17 PM	Standard	Aug 16, 2023, 4:34:17 PM	Not public	Download
part-00001-3079ec3f-c161-4cdd-...	application/octet-stream	Aug 16, 2023, 4:36:16 PM	Standard	Aug 16, 2023, 4:36:16 PM	Not public	Download
part-00001-4ef017bc-917d-42f7-b...	application/octet-stream	Aug 16, 2023, 4:34:18 PM	Standard	Aug 16, 2023, 4:34:18 PM	Not public	Download
part-00002-3079ec3f-c161-4cdd-...	application/octet-stream	Aug 16, 2023, 4:36:16 PM	Standard	Aug 16, 2023, 4:36:16 PM	Not public	Download
part-00002-4ef017bc-917d-42f7-b...	application/octet-stream	Aug 16, 2023, 4:34:20 PM	Standard	Aug 16, 2023, 4:34:20 PM	Not public	Download
part-00003-3079ec3f-c161-4cdd-...	application/octet-stream	Aug 16, 2023, 4:36:19 PM	Standard	Aug 16, 2023, 4:36:19 PM	Not public	Download
part-00003-4ef017bc-917d-42f7-b...	application/octet-stream	Aug 16, 2023, 4:34:22 PM	Standard	Aug 16, 2023, 4:34:22 PM	Not public	Download
part-00004-3079ec3f-c161-4cdd-...	application/octet-stream	Aug 16, 2023, 4:36:19 PM	Standard	Aug 16, 2023, 4:36:19 PM	Not public	Download
part-00004-4ef017bc-917d-42f7-b...	application/octet-stream	Aug 16, 2023, 4:34:22 PM	Standard	Aug 16, 2023, 4:34:22 PM	Not public	Download
part-00005-3079ec3f-c161-4cdd-...	application/octet-stream	Aug 16, 2023, 4:36:21 PM	Standard	Aug 16, 2023, 4:36:21 PM	Not public	Download
part-00005-4ef017bc-917d-42f7-b...	application/octet-stream	Aug 16, 2023, 4:34:27 PM	Standard	Aug 16, 2023, 4:34:27 PM	Not public	Download

Step 01. Prepare the MapReduce Java Application. In this case using Eclipse IDE.

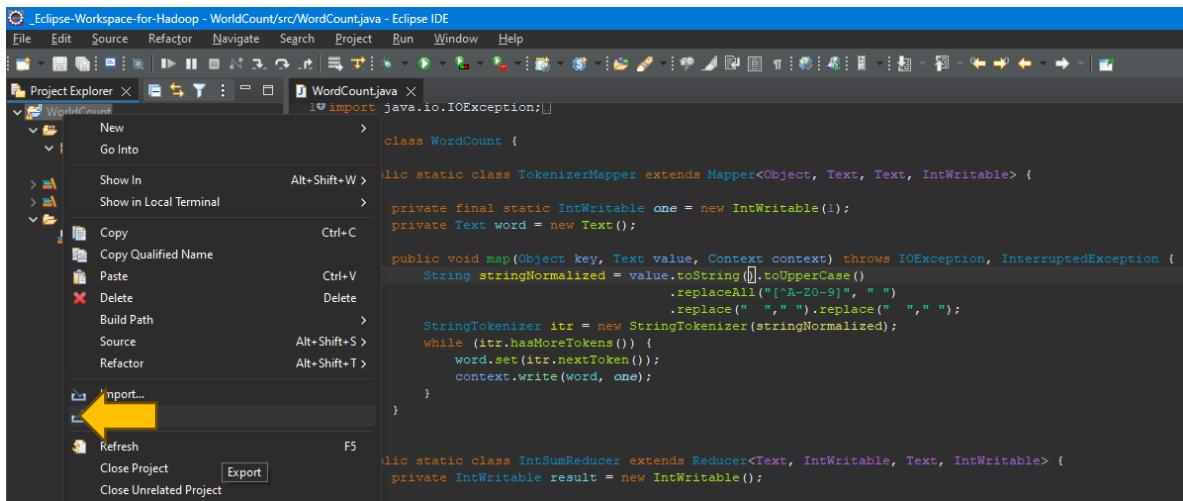


```

_Eclipse-Workspace-for-Hadoop - WordCount/src/WordCount.java - Eclipse IDE
File Edit Source Refactor Navigate Search Project Run Window Help
Project Explorer X WordCount.java X
WordCount
src
  (default package)
    WordCount.java
JRE System Library [jre7]
Referenced Libraries
Libraries
  hadoop-client-api-3.3.0.jar
WordCount.java
1 import java.io.IOException;
2
3 public class WordCount {
4
5     public static class TokenizerMapper extends Mapper<Object, Text, Text, IntWritable> {
6
7         private final static IntWritable one = new IntWritable(1);
8         private Text word = new Text();
9
10        public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
11            String stringNormalized = value.toString().toUpperCase()
12                .replaceAll("[A-Z0-9]", " ")
13                .replace(" ", " ").replace(" ", " ");
14            StringTokenizer itr = new StringTokenizer(stringNormalized);
15            while (itr.hasMoreTokens()) {
16                word.set(itr.nextToken());
17                context.write(word, one);
18            }
19        }
20
21        public static class IntSumReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
22            private IntWritable result = new IntWritable();
23
24            public void reduce(Text key, Iterable<IntWritable> values, Context context)
25                throws IOException, InterruptedException {
26                int sum = 0;
27                for (IntWritable value: values) { sum += value.get(); }
28                result.set(sum);
29                context.write(key, result);
30            }
31        }
32
33        public static void main(String[] args) throws Exception {
34            Configuration conf = new Configuration();
35            Job job = Job.getInstance(conf, "word count");
36            job.setJarByClass(WordCount.class);
37            job.setMapperClass(TokenizerMapper.class);
38            job.setCombinerClass(IntSumReducer.class);
39            job.setReducerClass(IntSumReducer.class);
40            job.setOutputKeyClass(Text.class);
41            job.setOutputValueClass(IntWritable.class);
42            FileInputFormat.addInputPath(job, new Path(args[0]));
43            FileOutputFormat.setOutputPath(job, new Path(args[1]));
44            System.exit(job.waitForCompletion(true) ? 0 : 1);
45        }
46    }
}

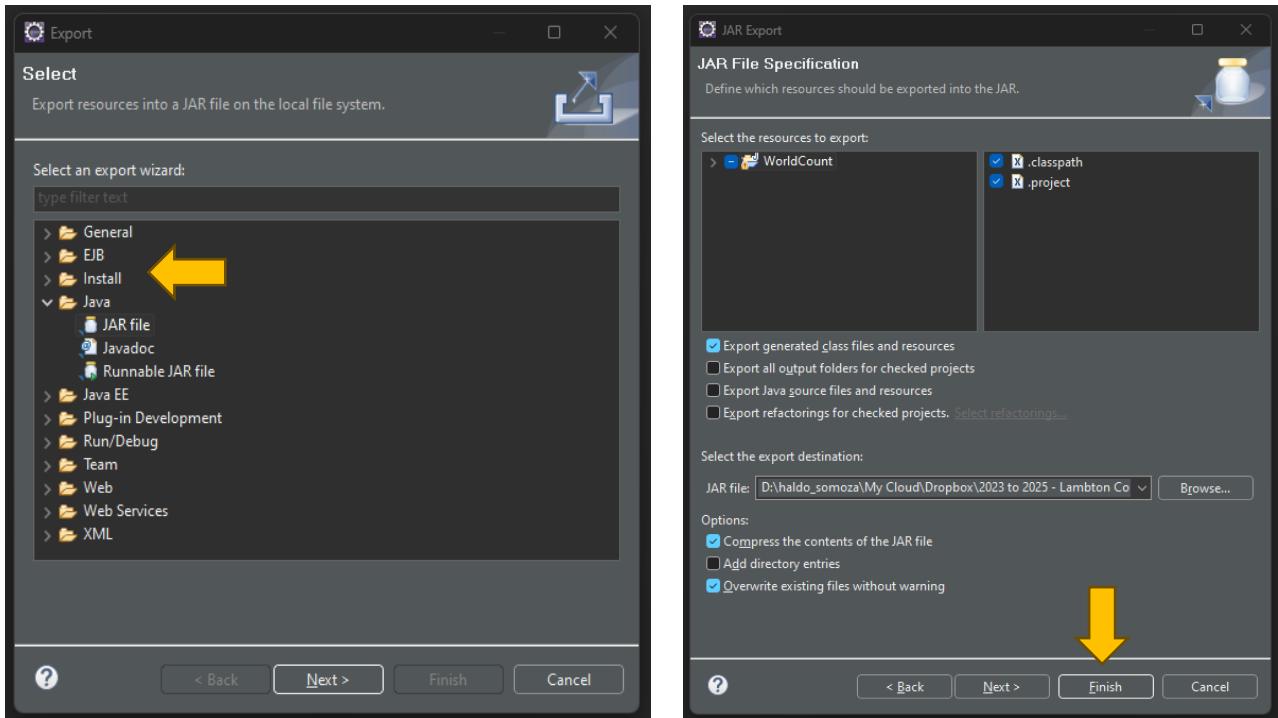
```

Right click over the project, choose Export

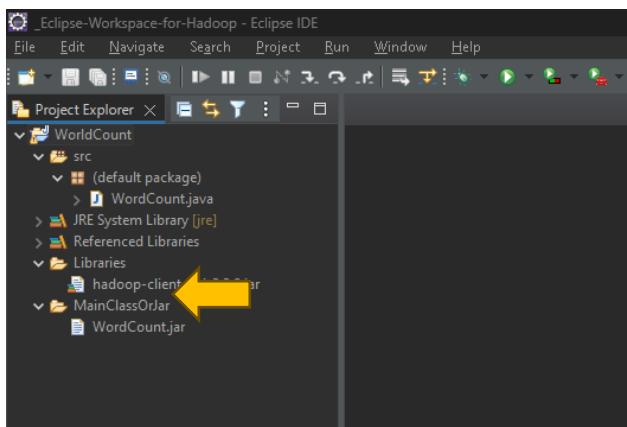


Choose Java -> JAR file. Press Next button.

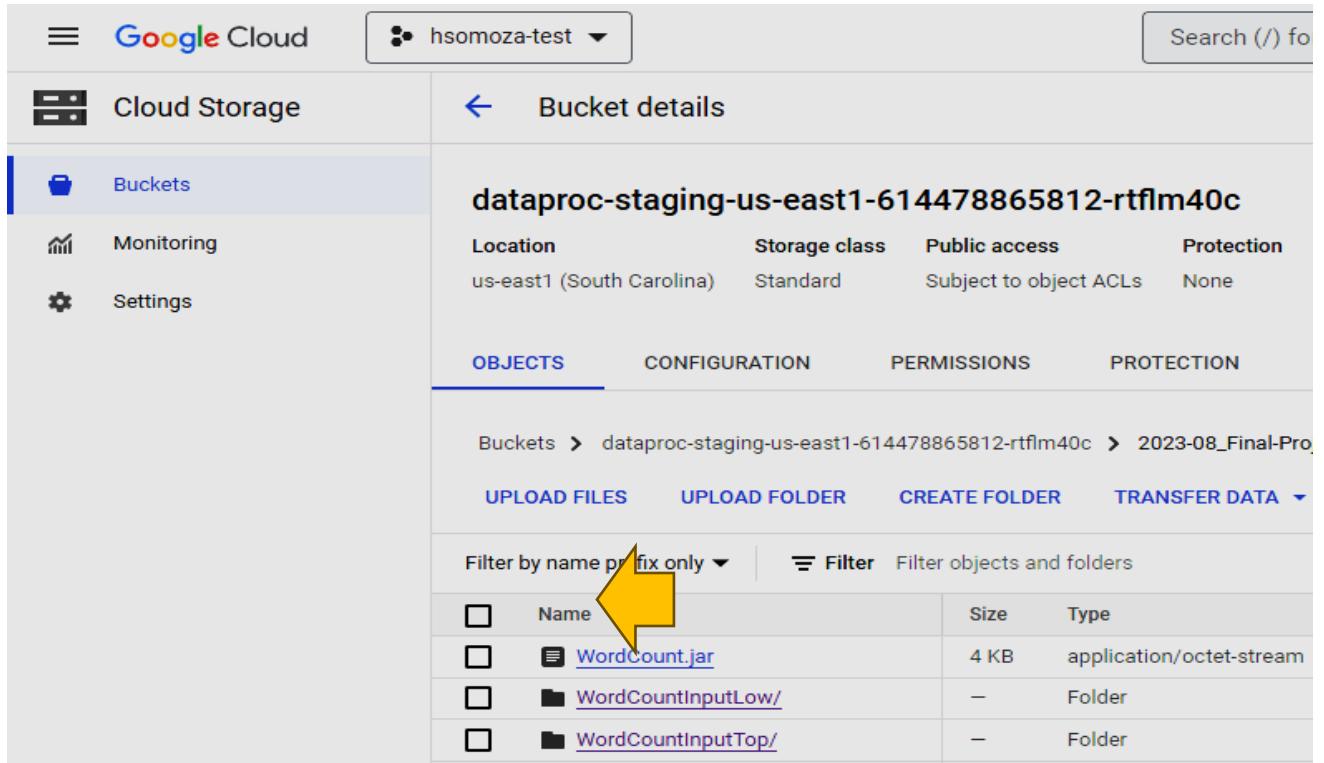
Select options like you see in the image and define one export destination folder and filename for the resulting JAR file. Then click on Finish.



Verify the generated JAR file.



Step 02: Upload the JAR file and the Input Data folders to a GCP Bucket.



The screenshot shows the Google Cloud Storage interface for the bucket 'hsomoza-test'. The left sidebar has 'Cloud Storage' selected. The main area displays the bucket details for 'dataproc-staging-us-east1-614478865812-rtflm40c'. The 'OBJECTS' tab is active, showing a list of uploaded files and folders:

Name	Size	Type
WordCount.jar	4 KB	application/octet-stream
WordCountInputLow/	—	Folder
WordCountInputTop/	—	Folder

A yellow arrow points to the 'Name' column header in the table. The URL in the browser's address bar is: <https://console.cloud.google.com/storage/buckets/dataproc-staging-us-east1-614478865812-rtflm40c/object?tab=objects>.

In this case upload two Input Data folders, one to process data with low ratings and another to process data with top ratings. In next image the WordCountInputLow folder with data generated by PySpark.

Bucket details

dataproj-staging-us-east1-614478865812-rtflm40c

Location	Storage class	Public access	Protection
us-east1 (South Carolina)	Standard	Subject to object ACLs	None

OBJECTS

Name	Size	Type	Created	Storage class	Last modified	Public access
yelp_results_mapr_part-00000-1b...	62.4 MB	text/csv	Aug 16, 2023, 5:29:38 PM	Standard	Aug 16, 2023, 5:29:38 PM	Not public
yelp_results_mapr_part-00001-1b...	57.1 MB	text/csv	Aug 16, 2023, 5:29:45 PM	Standard	Aug 16, 2023, 5:29:45 PM	Not public
yelp_results_mapr_part-00002-1b...	63 MB	text/csv	Aug 16, 2023, 5:30:51 PM	Standard	Aug 16, 2023, 5:30:51 PM	Not public
yelp_results_mapr_part-00003-1b...	58.4 MB	text/csv	Aug 16, 2023, 5:31:02 PM	Standard	Aug 16, 2023, 5:31:02 PM	Not public
yelp_results_mapr_part-00004-1b...	63.7 MB	text/csv	Aug 16, 2023, 5:31:11 PM	Standard	Aug 16, 2023, 5:31:11 PM	Not public
yelp_results_mapr_part-00005-1b...	59.6 MB	text/csv	Aug 16, 2023, 5:31:43 PM	Standard	Aug 16, 2023, 5:31:43 PM	Not public
yelp_results_mapr_part-00006-1b...	61 MB	text/csv	Aug 16, 2023, 5:31:55 PM	Standard	Aug 16, 2023, 5:31:55 PM	Not public
yelp_results_mapr_part-00007-1b...	60 MB	text/csv	Aug 16, 2023, 5:32:04 PM	Standard	Aug 16, 2023, 5:32:04 PM	Not public
yelp_results_mapr_part-00008-1b...	61.8 MB	text/csv	Aug 16, 2023, 5:32:39 PM	Standard	Aug 16, 2023, 5:32:39 PM	Not public
yelp_results_mapr_part-00009-1b...	56.5 MB	text/csv	Aug 16, 2023, 5:32:44 PM	Standard	Aug 16, 2023, 5:32:44 PM	Not public
yelp_results_mapr_part-00010-1b...	62.6 MB	text/csv	Aug 16, 2023, 5:33:00 PM	Standard	Aug 16, 2023, 5:33:00 PM	Not public
yelp_results_mapr_part-00011-1b...	58.7 MB	text/csv	Aug 16, 2023, 5:33:33 PM	Standard	Aug 16, 2023, 5:33:33 PM	Not public
Marketplace						
Release Notes						

...
 Marketplace
 Release Notes

Step 03: Creating and Executing one Hadoop Job for each two folders.

In Job definition window, choose Hadoop Job type, write the path of jar in Main class or jar box,

and such Arguments send (a) the Class Name, (b) the Input Data folder, and (c) the Output Data folder.

Google Cloud hsomoza-test Search (/) for resources, docs, products, and more

Dataproc

Submit a job

Jobs on Clusters

- Clusters
- Jobs**
- Workflows
- Autoscaling policies

Serverless

- Batches

Metastore Services

- Metastore
- Federation

Utilities

- Component exchange
- Workbench

Release Notes

Job ID * job-ddc5f96c-WordCountInputLow

Region * us-east1

Specifies the Cloud Dataproc regional service, which determines what clusters are available.

Cluster * cluster-2702

Job type * Hadoop

Main class or jar * gs://dataproc-staging-us-east1-614478865812-rtflm40c/2023-08_Final-Project-BDM-1

The fully qualified name of a class in a provided or standard jar file, for example, com.example.wordcount, or a provided jar file to use the main class of that jar file

Jar files

Jar files are included in the CLASSPATH. Can be a GCS file with the gs:// prefix, an HDFS file on the cluster with the hdfs:// prefix, or a local file on the cluster with the file:// prefix.

Archive files

Archive files are extracted in the Spark working directory. Can be a GCS file with the gs:// prefix, an HDFS file on the cluster with the hdfs:// prefix, or a local file on the cluster with the file:// prefix. Supported file types: jar, .tar, .tar.gz, .tgz, .zip.

Arguments

- WordCount
- gs://dataproc-staging-us-east1-614478865812-rtflm40c/2023-08_Final-Project-BDM-1024/WordCountInputLow/
- gs://dataproc-staging-us-east1-614478865812-rtflm40c/2023-08_Final-Project-BDM-1024/WordCountOutputLow/

Additional arguments to pass to the main class. Press Return after each argument.

Max restarts per hour

Leave blank if you don't want to allow automatic restarts on job failure. [Learn more](#)

Properties

+ ADD PROPERTY

Labels

+ ADD LABEL

SUBMIT CANCEL

EQUIVALENT REST

The jobs finished successfully, showing the following information and log, similarly in both jobs.



The log resulting:

```
2023-08-16 22:29:28,765 INFO client.DefaultNoHARMF failoverProxyProvider: Connecting to ResourceManager at
cluster-2702-m.us-east1-b.c.hsomoza-test.internal./10.142.0.4:8032

2023-08-16 22:29:29,066 INFO client.AHSProxy: Connecting to Application History server at cluster-2702-m.us-
east1-b.c.hsomoza-test.internal./10.142.0.4:10200

Aug 16, 2023 10:29:29 PM com.google.cloud.hadoop.fs.gcs.GhfsStorageStatistics updateStats
WARNING: Detected potential high latency for operation op_get_file_status. latencyMs=361;
previousMaxLatencyMs=0; operationCount=1; context=gs://dataproc-staging-us-east1-614478865812-rtflm40c/2023-
08_Final-Project-BDM-1024/WordCountOutputLow

2023-08-16 22:29:30,185 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not
performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.

2023-08-16 22:29:30,205 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path:
/tmp/hadoop-yarn/staging/root/.staging/job_1692216534885_0007

2023-08-16 22:29:30,684 INFO input.FileInputFormat: Total input files to process : 30

2023-08-16 22:29:30,795 INFO mapreduce.JobSubmitter: number of splits:30

2023-08-16 22:29:31,072 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1692216534885_0007

2023-08-16 22:29:31,072 INFO mapreduce.JobSubmitter: Executing with tokens: []

2023-08-16 22:29:31,317 INFO conf.Configuration: resource-types.xml not found

2023-08-16 22:29:31,317 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.

2023-08-16 22:29:31,597 INFO impl.YarnClientImpl: Submitted application application_1692216534885_0007

2023-08-16 22:29:31,630 INFO mapreduce.Job: The url to track the job: http://cluster-2702-m.us-east1-
b.c.hsomoza-test.internal.:8088/proxy/application_1692216534885_0007/

2023-08-16 22:29:31,631 INFO mapreduce.Job: Running job: job_1692216534885_0007

2023-08-16 22:29:42,740 INFO mapreduce.Job: Job job_1692216534885_0007 running in uber mode : false

2023-08-16 22:29:42,741 INFO mapreduce.Job: map 0% reduce 0%

2023-08-16 22:30:00,899 INFO mapreduce.Job: map 2% reduce 0%

2023-08-16 22:30:05,926 INFO mapreduce.Job: map 4% reduce 0%

2023-08-16 22:30:10,977 INFO mapreduce.Job: map 5% reduce 0%

2023-08-16 22:30:11,983 INFO mapreduce.Job: map 7% reduce 0%
```

```
2023-08-16 22:30:29,082 INFO mapreduce.Job: map 9% reduce 0%
2023-08-16 22:30:30,087 INFO mapreduce.Job: map 10% reduce 0%
2023-08-16 22:30:33,102 INFO mapreduce.Job: map 11% reduce 0%
2023-08-16 22:30:34,115 INFO mapreduce.Job: map 12% reduce 0%
2023-08-16 22:30:39,145 INFO mapreduce.Job: map 13% reduce 0%
2023-08-16 22:30:54,226 INFO mapreduce.Job: map 14% reduce 0%
2023-08-16 22:30:56,234 INFO mapreduce.Job: map 16% reduce 0%
2023-08-16 22:30:57,239 INFO mapreduce.Job: map 17% reduce 0%
2023-08-16 22:31:00,255 INFO mapreduce.Job: map 18% reduce 0%
2023-08-16 22:31:06,284 INFO mapreduce.Job: map 20% reduce 0%
2023-08-16 22:31:18,350 INFO mapreduce.Job: map 21% reduce 0%
2023-08-16 22:31:23,377 INFO mapreduce.Job: map 24% reduce 0%
2023-08-16 22:31:24,382 INFO mapreduce.Job: map 25% reduce 0%
2023-08-16 22:31:32,420 INFO mapreduce.Job: map 27% reduce 0%
2023-08-16 22:31:44,483 INFO mapreduce.Job: map 28% reduce 0%
2023-08-16 22:31:46,492 INFO mapreduce.Job: map 29% reduce 0%
2023-08-16 22:31:50,510 INFO mapreduce.Job: map 31% reduce 0%
2023-08-16 22:31:52,517 INFO mapreduce.Job: map 32% reduce 0%
2023-08-16 22:31:58,538 INFO mapreduce.Job: map 33% reduce 0%
2023-08-16 22:31:59,542 INFO mapreduce.Job: map 34% reduce 0%
2023-08-16 22:32:11,599 INFO mapreduce.Job: map 36% reduce 0%
2023-08-16 22:32:12,605 INFO mapreduce.Job: map 37% reduce 0%
2023-08-16 22:32:17,627 INFO mapreduce.Job: map 39% reduce 0%
2023-08-16 22:32:25,658 INFO mapreduce.Job: map 40% reduce 0%
2023-08-16 22:32:31,681 INFO mapreduce.Job: map 41% reduce 0%
2023-08-16 22:32:33,689 INFO mapreduce.Job: map 42% reduce 0%
2023-08-16 22:32:37,705 INFO mapreduce.Job: map 43% reduce 0%
2023-08-16 22:32:43,730 INFO mapreduce.Job: map 46% reduce 0%
```

```
2023-08-16 22:32:48,750 INFO mapreduce.Job: map 47% reduce 0%
2023-08-16 22:32:55,780 INFO mapreduce.Job: map 48% reduce 0%
2023-08-16 22:32:58,796 INFO mapreduce.Job: map 49% reduce 0%
2023-08-16 22:32:59,799 INFO mapreduce.Job: map 50% reduce 0%
2023-08-16 22:33:06,826 INFO mapreduce.Job: map 52% reduce 0%
2023-08-16 22:33:10,840 INFO mapreduce.Job: map 53% reduce 0%
2023-08-16 22:33:18,874 INFO mapreduce.Job: map 55% reduce 0%
2023-08-16 22:33:20,880 INFO mapreduce.Job: map 56% reduce 0%
2023-08-16 22:33:24,893 INFO mapreduce.Job: map 57% reduce 0%
2023-08-16 22:33:26,901 INFO mapreduce.Job: map 58% reduce 0%
2023-08-16 22:33:27,905 INFO mapreduce.Job: map 60% reduce 0%
2023-08-16 22:33:30,917 INFO mapreduce.Job: map 61% reduce 0%
2023-08-16 22:33:36,939 INFO mapreduce.Job: map 62% reduce 0%
2023-08-16 22:33:37,942 INFO mapreduce.Job: map 63% reduce 0%
2023-08-16 22:33:47,974 INFO mapreduce.Job: map 66% reduce 0%
2023-08-16 22:33:50,984 INFO mapreduce.Job: map 67% reduce 0%
2023-08-16 22:33:57,002 INFO mapreduce.Job: map 68% reduce 0%
2023-08-16 22:34:00,018 INFO mapreduce.Job: map 69% reduce 0%
2023-08-16 22:34:03,028 INFO mapreduce.Job: map 70% reduce 0%
2023-08-16 22:34:06,037 INFO mapreduce.Job: map 71% reduce 0%
2023-08-16 22:34:08,043 INFO mapreduce.Job: map 73% reduce 0%
2023-08-16 22:34:11,055 INFO mapreduce.Job: map 74% reduce 0%
2023-08-16 22:34:15,072 INFO mapreduce.Job: map 76% reduce 0%
2023-08-16 22:34:16,076 INFO mapreduce.Job: map 77% reduce 0%
2023-08-16 22:34:29,120 INFO mapreduce.Job: map 79% reduce 0%
2023-08-16 22:34:32,129 INFO mapreduce.Job: map 80% reduce 0%
2023-08-16 22:34:36,141 INFO mapreduce.Job: map 81% reduce 0%
2023-08-16 22:34:38,147 INFO mapreduce.Job: map 83% reduce 0%
```

```
2023-08-16 22:34:42,160 INFO mapreduce.Job: map 84% reduce 0%
2023-08-16 22:34:49,180 INFO mapreduce.Job: map 87% reduce 0%
2023-08-16 22:34:51,188 INFO mapreduce.Job: map 88% reduce 0%
2023-08-16 22:34:53,196 INFO mapreduce.Job: map 90% reduce 0%
2023-08-16 22:35:08,254 INFO mapreduce.Job: map 92% reduce 0%
2023-08-16 22:35:10,260 INFO mapreduce.Job: map 93% reduce 0%
2023-08-16 22:35:15,274 INFO mapreduce.Job: map 97% reduce 0%
2023-08-16 22:35:21,291 INFO mapreduce.Job: map 98% reduce 0%
2023-08-16 22:35:26,305 INFO mapreduce.Job: map 100% reduce 0%
2023-08-16 22:35:39,349 INFO mapreduce.Job: map 100% reduce 33%
2023-08-16 22:35:45,366 INFO mapreduce.Job: map 100% reduce 67%
2023-08-16 22:35:46,369 INFO mapreduce.Job: map 100% reduce 100%
2023-08-16 22:35:47,377 INFO mapreduce.Job: Job job_1692216534885_0007 completed successfully
2023-08-16 22:35:47,506 INFO mapreduce.Job: Counters: 61
```

File System Counters

```
FILE: Number of bytes read=38735724
```

```
FILE: Number of bytes written=71646869
```

```
FILE: Number of read operations=0
```

```
FILE: Number of large read operations=0
```

```
FILE: Number of write operations=0
```

```
GS: Number of bytes read=1538677788
```

```
GS: Number of bytes written=2298801
```

```
GS: Number of read operations=375670
```

```
GS: Number of large read operations=0
```

```
GS: Number of write operations=835544
```

```
HDFS: Number of bytes read=7260
```

```
HDFS: Number of bytes written=0
```

```
HDFS: Number of read operations=30
```

HDFS: Number of large read operations=0

HDFS: Number of write operations=0

HDFS: Number of bytes read erasure-coded=0

Job Counters

Killed map tasks=1

Killed reduce tasks=1

Launched map tasks=30

Launched reduce tasks=4

Rack-local map tasks=30

Total time spent by all maps in occupied slots (ms)=3221969339

Total time spent by all reduces in occupied slots (ms)=155028316

Total time spent by all map tasks (ms)=983207

Total time spent by all reduce tasks (ms)=47308

Total vcore-milliseconds taken by all map tasks=983207

Total vcore-milliseconds taken by all reduce tasks=47308

Total megabyte-milliseconds taken by all map tasks=3221969339

Total megabyte-milliseconds taken by all reduce tasks=155028316

Map-Reduce Framework

Map input records=4326178

Map output records=288118351

Map output bytes=2640412774

Map output materialized bytes=23539632

Input split bytes=7260

Combine input records=288118351

Combine output records=1681325

Reduce input groups=208886

Reduce shuffle bytes=23539632

Reduce input records=1681325

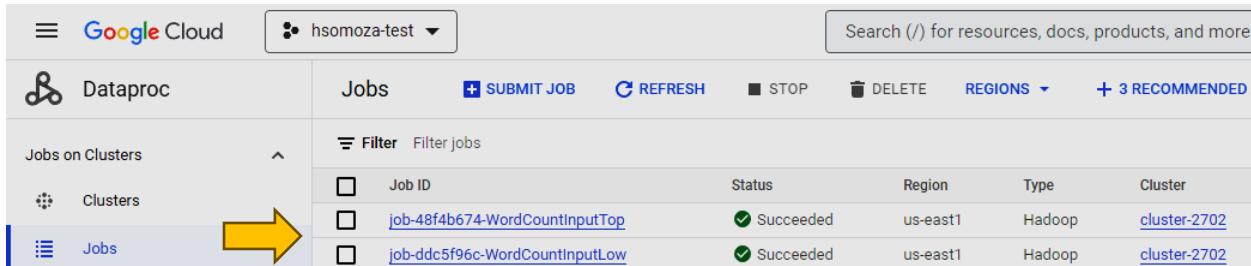
```
Reduce output records=208886
Spilled Records=4440584
Shuffled Maps =90
Failed Shuffles=0
Merged Map outputs=90
GC time elapsed (ms)=11007
CPU time spent (ms)=851400
Physical memory (bytes) snapshot=23218028544
Virtual memory (bytes) snapshot=156121165824
Total committed heap usage (bytes)=18877513728
Peak Map Physical memory (bytes)=793665536
Peak Map Virtual memory (bytes)=4743344128
Peak Reduce Physical memory (bytes)=502362112
Peak Reduce Virtual memory (bytes)=4730691584

Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0

File Input Format Counters
Bytes Read=1538677788

File Output Format Counters
Bytes Written=2298801
```

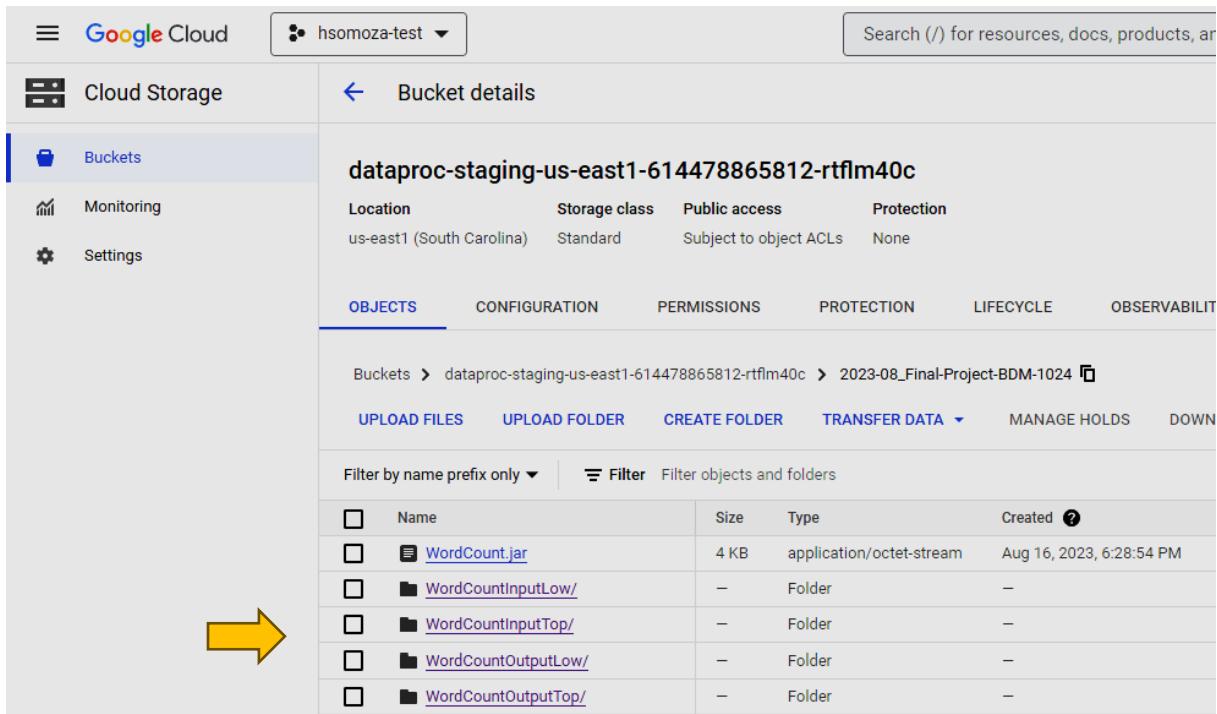
Both Hadoop jobs finished successfully.



The screenshot shows the Google Cloud DataProc Jobs page. The left sidebar has 'Dataproc' selected under 'Jobs'. A yellow arrow points from the 'Jobs' section in the sidebar to the list of jobs on the right. The list shows two entries:

Job ID	Status	Region	Type	Cluster
job-48f4b674-WordCountInputTop	Succeeded	us-east1	Hadoop	cluster-2702
job-ddc5f96c-WordCountInputLow	Succeeded	us-east1	Hadoop	cluster-2702

Step 04: Identifying and Importing the resulting files.



The screenshot shows the Google Cloud Storage Bucket details page for 'dataproc-staging-us-east1-614478865812-rtflm40c'. The left sidebar has 'Buckets' selected. A yellow arrow points from the 'Buckets' section in the sidebar to the list of objects below. The list shows the following files:

Name	Size	Type	Created
WordCount.jar	4 KB	application/octet-stream	Aug 16, 2023, 6:28:54 PM
WordCountInputLow/	—	Folder	—
WordCountInputTop/	—	Folder	—
WordCountOutputLow/	—	Folder	—
WordCountOutputTop/	—	Folder	—

Here the result for one job splatted in three parts.

Google Cloud Bucket details for **dataproc-staging-us-east1-614478865812-rtflm40c**

Name	Size	Type	Created
_SUCCESS	0 B	application/octet-stream	Aug 16, 2023, 6:35:44 PM
part-r-00000	743.7 KB	application/octet-stream	Aug 16, 2023, 6:35:37 PM
part-r-00001	748.4 KB	application/octet-stream	Aug 16, 2023, 6:35:43 PM
part-r-00002	752.8 KB	application/octet-stream	Aug 16, 2023, 6:35:44 PM

Step 05: Loading the data into one Excel file to analyze.

The downloaded data was joined in one only CSV file and then loaded into one Excel file.

Eclipse-Workspace-for-Hadoop - WorldCount/WordCountOutputLow/2023-08_Final-Project-BDM-1024_WordCountOutputLow_part-r-00000 - Eclipse IDE

Project Explorer

- WorldCount
 - src
 - JRE System Library [jre]
 - Referenced Libraries
 - Libraries
 - MainClassOrJar
 - WordCountInput
 - WordCountInputLow
 - WordCountInputTop
 - WordCountOutput
 - WordCountOutputLow
 - _SUCCESS
 - 2023-08_Final-Project-BDM-1024_WordCountOutputLow_part-r-00000
 - 2023-08_Final-Project-BDM-1024_WordCountOutputLow_part-r-00001
 - 2023-08_Final-Project-BDM-1024_WordCountOutputLow_part-r-00002
 - 2023-08_Final-Project-BDM-1024_WordCountOutputLow.csv
 - 2023-08_Final-Project-BDM-1024_WordCountOutputLow.xlsx
 - WordCountOutputTop
 - _SUCCESS
 - 2023-08_Final-Project-BDM-1024_WordCountOutputTop_part-r-00000
 - 2023-08_Final-Project-BDM-1024_WordCountOutputTop_part-r-00001
 - 2023-08_Final-Project-BDM-1024_WordCountOutputTop_part-r-00002
 - 2023-08_Final-Project-BDM-1024_WordCountOutputTop.csv
 - 2023-08_Final-Project-BDM-1024_WordCountOutputTop.xlsx

2023-08_Final-Project-BDM-1024_WordCountOutputLow_part-r-00000

10000 BRONZERO	1
10001 BRONZER	2
10002 BRONZIE	2
10003 BRONZINOS	1
10004 BROOCH	3
10005 BROOD	26
10006 BRODMORE	1
10007 BROODY	6
10008 BROOKIES	3
10009 BROOKLAWN	32
10010 BROOKLY	3
10011 BROOKLYN	3893
10012 BROOKLYNN	7
10013 BROOKSVILLE	59
10014 BROOKVILLE	14
10015 BROOM	313
10016 BROOMALL	226
10017 BROOMHALL	1
10018 BROOMS	30
10019 BROUGHT	7
10020 BROPHEY	1
10021 BRORIPPLE	3
10022 BROSEPH	6
10023 BROSNAHAN	5
10024 BROSO	2
10025 BROSTEP	3
10026 BROTH	9

	Column1	Column2
1	THE	14,979,940
2	AND	10,876,489
3	I	7,542,451
4	A	7,075,533
5	TO	6,147,001
6	WAS	6,128,502
7	IT	4,236,889
8	IS	3,697,106
9	OF	3,552,663
10	FOR	3,431,234
11	WE	3,188,934
12	IN	2,875,425
13	FOOD	2,616,651
14	THIS	2,430,202
15	MY	2,349,012
16	BUT	2,300,513
17	WITH	2,221,532
18	THEY	2,201,120
19	HAD	2,079,112
20	GOOD	1,952,639
21	ON	1,941,847
22	GREAT	1,938,764
23	THAT	1,932,030
24	WERE	1,908,635
25	T	1,884,233
26	YOU	1,856,532
27	NOT	1,756,340
28	PLACE	1,718,896

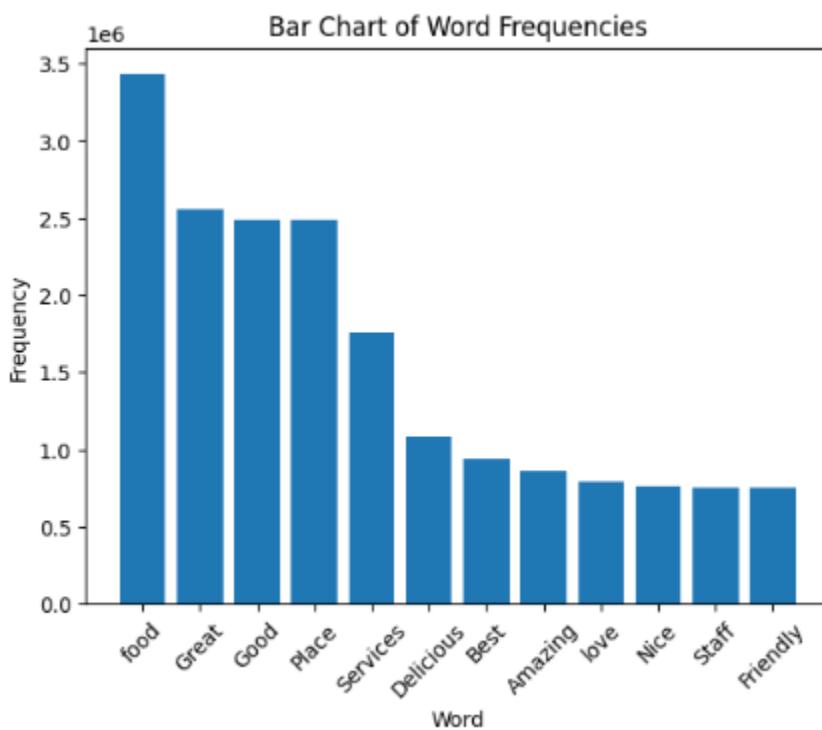
	Column1	Column2
1	THE	19,780,139
2	AND	14,473,770
3	I	9,963,887
4	A	9,085,506
5	WAS	8,014,275
6	TO	7,702,173
7	IT	5,531,430
8	IS	5,009,774
9	OF	4,781,837
10	FOR	4,450,545
11	WE	3,985,476
12	IN	3,811,599
13	FOOD	3,429,798
14	THIS	3,190,459
15	MY	3,032,301
16	BUT	2,986,878
17	WITH	2,971,042
18	HAD	2,773,730
19	THEY	2,685,517
20	GREAT	2,552,985
21	GOOD	2,484,223
22	YOU	2,477,185
23	THAT	2,438,569
24	WERE	2,426,199
25	ON	2,395,694
26	PLACE	2,372,759
27	T	2,312,697
28	NOT	2,138,115

By ran a MapReduce job to count the words of good_rate and bad_rate reviews we just understand that there are lots of common words that we need to clarify them by an advance Machine Learning algorithms to find out the specific result, so for now we just manually went trough all the result and remove the non-relative words that may not have a value in this result.

For example we found out for the good rated reviews people mostly used these worded

A	B	C	D
Column1	Column2		
14 FOOD	3429798		
20 THEY	2685517		
21 GREAT	2552985		
22 GOOD	2484223		
27 PLACE	2372759		
35 SERVICE	1758807		
50 DELICIOUS	1087390		
56 BEST	943263		
61 AMAZING	857580		
69 LOVE	793568		
70 GOT	779726		
71 VE	768092		
72 NICE	767495		
73 BEEN	762117		
74 DEFINITELY	751551		
75 STAFF	751292		
76 FRIENDLY	751161		
77 CHICKEN	741094		
78 UP	736032		

Food, Place, Service, Staff So it just shows us people don't evaluate a restaurant just based on the food, they also pay attention to services the place and the manner of the staff as well.



The most common word in the bar graph is "food," with a frequency of 3429798. This suggests that food is the most important factor for customers when rating a restaurant. It also provides some insights into what customers are looking for in a restaurant. For example, the words "place," "services," and "staff" suggest that customers also value the atmosphere and service of a restaurant. The words "delicious," "best," and "amazing" suggest that customers are looking for a high-quality dining experience. Keep in mind that this is just top written words on good-rated reviews.

Conclusions and Recommendations

It was identified that the city with the largest number of restaurants is Philadelphia with a total of 5,852 businesses in this category, with an approximate average of 3.56 stars. It was also found that 7 of the 10 cities with the largest number of restaurants have more than 2,000 business units of this type. It is recommended that the company make a strategic approach to increase the number of restaurants that want to pay for advertising in these cities, and likewise make a strategic plan to increase user participation in them. This benefits both financially and in growth and strengthening of the YELP brand.

1,553 restaurants were found with a high 5-star rating from their users, which represents about the 3% of the population. It is recommended that the company take advantage of this information to establish a closer collaboration with these featured restaurants. They could analyze what is making the difference in these restaurants and identify a strategy to help other business to success. Collaborations must be organized based on special and strategic marketing programs. Yelp will also benefit as recommending these highly rated restaurants on its website can increase user satisfaction and loyalty and strengthen its position as a trusted platform for finding the best restaurants.

56079 restaurant reviews were found, with a star rating equal to or less than 2. 35% of these reviews are related to the issues “Customer services” or “Clean” or “Wait”. MacDonald’s, Taco Bell, and Wendy’s are the most repetitive restaurants with rates under 2. It is recommended to develop strategic actions to correct or minimize problems in the 58, 45, and 296 restaurants rated with stars less than or

equal to 2 and more than 10 reviews on issues related to "customer service", "cleanliness" and "waits" respectively. This may include setting quality guidelines for Yelp-listed restaurants, providing additional training for customer service staff, and encouraging improved cleaning and customer service standards at partner establishments. The implementation of the recommendations at this point would mean achieving and strengthening the reputation of businesses on the platform and increasing customer satisfaction, which translates customer loyalty to the platform, reducing user turnover and lowering the costs associated with the acquisition of new users, which benefits the profitability of the company and could also lead to an increase in advertising exposure and therefore advertising revenue.

The 10 users with the largest number of fans were found, among these the one with the most followers has a total of 12497, approximately 29% more than the user with the second most fans. All these users have more than 2,000 followers, a significant number of fans. It is recommended to establish solid relationships with the most influential users and develop marketing plans that include creating advertising guidelines with added value for interested businesses, this would include recommending businesses based on valuable content, and interacting with other users, among others. This would allow YELP to charge special rates for this type of advertising due to the reach and reputation these influential users have, creating an additional revenue stream for the company.

The 100 users with the largest number of reviews were obtained, the first of which is called Fox and has 17,473 reviews. Yelp is encouraged to design different rewards programs for users who make a high number of reviews. This can include benefits such as exclusive discounts, special promotions, or even points that can be redeemed for additional prizes or incentives. They could also create referral programs,

among others. All this to improve loyalty to the platform, which can in turn generate an expansion of its user base and generate greater economic opportunities.

To conclude this analytical solution focus on customer rating enhancement that creates a mutually beneficial environment, improving Yelp's reputation and generating greater revenue opportunities for businesses and for Yelp as well.

Roles of Group Members

TASKS	RESPONSIBLES
Design the postgresql database (model diagram)	Eduardo Williams / Adriana Penaranda
Design a python program to transform the information from json to Postgresql database	Haldo Somoza / Nilesh Khurana
Create the database in the cloud	Haldo Somoza
Import the information to Postgresql in the cloud	Marzieh Mohammadi / Adriana Penaranda / Haldo Somza / Eduardo Williams
Cleaning and filtering. (Restaurants) and check the quality of the data	Ishika Sukhija / Carlos Rey
Design SQL queries	Carlos Rey / Luis Gutierrez / Ishika Sukhija
Making the Mid Point Presentation	Marzieh Mohammadi / Adriana Penaranda/ Ishika Sukhija
Spark Jobs	Eduardo Williams / Adriana Penaranda
Map Reduce	Haldo Somoza / Marizeh Mohammadi
Visualization	Carlos Rey / Ishika Sukhija / Luis Gutierrez
Final Report	Marizeh Mohammadi / Nilesh Khurana
Final Presentation	Adriana Penaranda

References

yelp. (2023). Retrieved from <https://www.yelp.com/dataset>