# RBDA Project Part 1: Data Ingestion

**Name:** Aygun Najafova     **NetID:** an4758     **Group:** 21

## Idea

As a group, we are working on analysing important data points, including NYC restaurant inspection data, NYPD arrest data, NYC hotel reviews, and MTA transit data. From the above data sources, my dataset is the NYC Inspection data, and I have completed the data ingestion phase for the same. It is explained in more detail in the next sections.

## Data source

I am using the `DOHMH New York City Restaurant Inspection Results` dataset available on the NYC OpenData website.[1] As shown on the dataset's webpage, it contains detailed information about restaurant inspections in New York City, including violations and assigned grades, collected over more than ten years.

According to the website, the dataset contains over 291,000 rows and 27 columns. Each row is a single restaurant inspection record, and the columns have information on restaurant details, location, scoring and some additional information. These columns are explained below:

- **CAMIS**: Its a unique restaurant identifier stored as a text field.

- **DBA**: It is the restaurant name, represented as text.

- **BORO**: It is the borough name, stored as a text value.

- **BUILDING**: It is the number of the building in text format.

- **STREET**: It is the Street name recorded as text.

- **ZIPCODE**: The ZIP code of the restaurant as text.

- **PHONE**: The phone number of the place stored as text.

- **CUISINE DESCRIPTION**: The category of cuisine represented as text.

- **INSPECTION DATE**: The inspection date stored as a timestamp.

- **ACTION**: The outcome of the inspection stored as text.

- **VIOLATION CODE**: The code to identify violation, stored as text.

- **VIOLATION DESCRIPTION**: The description of the violation, stored as text.

- **CRITICAL FLAG**: Severity of the violation, represented as text.

- **SCORE**: The score of the inspection stored as an integer.

- **GRADE**: Grades for the inspection (A, B, C) represented as text.

- **GRADE DATE**: The date the grade was issued, as a timestamp.

- **RECORD DATE**: The date the record was entered or updated as a timestamp.

- **INSPECTION TYPE**: The type of inspection stored as text.

- **Latitude**: The latitude coordinate stored as a decimal number.

- **Longitude**: The longitude coordinate stored as a decimal number.

- Other Metadata Fields like `Community Board`, `Council District`, `Census Tract`, `BIN`, `BBL`, `NTA`, and `Location` that represent geographic identifiers and are not important, so I will remove these columns in my preprocessing phase.

---

[1] Dataset available at: `https://data.cityofnewyork.us/Health/DOHMH-New-York-City-Restaurant-Inspection-Results/43nn-pn8j`
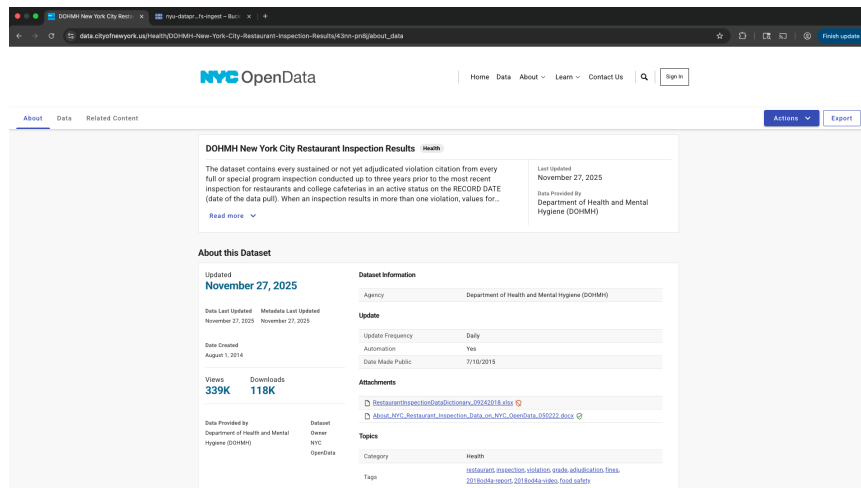
Figure 1: Data source on NYC OpenData

## Data analysis

To begin my analysis, I selected the first 100 rows of the dataset for a quick visual inspection. This was done using the following command:

```
head -n 101 NYC_Restaurent_Inscp_Data.csv > NYC_Restaurent_Inscp_Data_First100.csv
```

After opening the file in a spreadsheet, I observed several missing values in different columns. Some records had empty CUISINE DESCRIPTION fields, missing ZIPCODE, and even missing borough information. Also, important columns such as ZIPCODE and INSPECTION DATE were present which makes it suitable to index the data for future study of trends in the subsequent step of the project.

Additionally, I noticed that the dataset contains very old inspection dates, with some records dating back to the early 1900s. These entries are not useful to analyse recent trends. I am interested in recent and frequently monitored restaurants. Thus, I will only consider records from the year 2015 onwards.

Figure 2: Opening a subset of data in a spreadsheet

One more part of my analysis was actually counting the content and analysing the columns and the actual size of the data. I can see that it is a decently big dataset with over 292,256 rows and it is around 150 MB in size.



Figure 3: Basic analysis using shell commands

# Data Ingestion using MapReduce

Now that I have analysed the dataset, the next step is to write a MapReduce program to process the data and ingest it into the Hadoop system. This step is essential because the processed data will be used by Hive or Trino for further analysis in the next phase of the project, and it can also be loaded into Tableau for interactive visual analysis.

## Mapper

The mapper performs a number of operations, including checking for any missing values in the CSV and also transforming the missing values for some fields with valid data.

- **Input Processing**
  Each input in Mapper is the CSV row. It is received as a `Text` value. I am parsing using Apache Commons CSV with default comma-delimited, quote-aware settings.

- **Filtering and Validation**
  I am processing each row to perform validation on the schema of the data and only considering specific INPECTION DATE records.

  - The first entry of the Mapper is the CSV header row, which is skipped since it is of no use as it is already added at the Reducer side.
    I am also dropping the entire record from the dataset when:
    * The **SCORE** column is null or empty. This should be dropped since there is no suitable replacement for it.
    * The records with **INSPECTION DATE** before 2015 are dropped since I am only interested in the analysis of data from January 2015.
    * The **ZIPCODE** is an essential field, and thus I am dropping the records that have missing ZIPCODEs.
    * I am also dropping the records where **VIOLATION CODE** is null or empty.

- **Transformations**

  - All the values of the column fields in the records are trimmed and the nulls or missing values are replaced with empty strings.
  - For the case where **PHONE** is empty or missing, I am replacing it with 0 as a default number case.
  - The **SCORE** is a double value, which is converted to an Integer and stored in the format of Text string.
  - The **GRADE** is an important column but most of the records do not provide a grade or contain missing values. To fix this the missing grades, I calculated it using the official NYC Department of Health scoring guidelines using the below lookup.

$$
\text{GRADE} = \begin{cases} \text{A,} & 0 \leq \text{SCORE} \leq 13, \\ \text{B,} & 14 \leq \text{SCORE} \leq 27, \\ \text{C,} & \text{SCORE} \geq 28. \end{cases}
$$

  If a GRADE is already present, I will just remove the whitespace and process it. This ensures I always have a GRADE value in my records.

- **Column Selection**
  After processing all the columns and validating their schema, now I am only considering the following columns in my final output CSV to the reducer side.
  **Allowed Columns:**

  CAMIS, INSPECTION DATE, DBA, BORO, BUILDING, STREET, ZIPCODE, PHONE, CUISINE DESCRIPTION, ACTION, VIOLATION CODE, VIOLATION DESCRIPTION, CRITICAL FLAG, SCORE, GRADE, INSPECTION TYPE

  **Dropped Columns**

  GRADE DATE, RECORD DATE, Latitude, Longitude, Community Board, Council District, Census Tract, BIN, BBL, NTA, Location

- At the end, I am also formatting the Mapper output values by properly escaping the special characters like new line or intermediate ',' in column values.

  **Output Record Structure:**

  - **Key:** `NullWritable`
  - **Value:** CSV row string with all 16 selected columns.

- Additionally, I am also collecting a set of Mapper-level statistics to understand data quality, study some distributions, and validate the filtering logic using counters. It helps validate the correctness of the preprocessing step and provides insights into how the dataset is processed after cleaning. The counters are as follows:

  - **Input Counters:** These counters track the total number of raw input rows read by the mapper. It provides understanding of how many records were present in the input.
  - **Output Counters:** These counters count the number of records that were successfully cleaned and emitted as final processed output.
  - **Borough Counters:** These counters aggregate valid records by borough to provide some interesting insights into patterns in the dataset. For each borough, I am tracking:
    * Total number of valid inspection records per borough.
    * Number of records with a Grade A rating per borough.
    * Number of restaurants serving American cuisine per borough.

## Reducer

The reducer is responsible for producing the final cleaned output. The `NYCInspectionReducer` performs two tasks that is writing the header row and emitting all valid data records passed from the mapper.

### Writting Header row

During the `setup()` phase, the reducer writes the CSV header as the very first output line. This header lists the 16 selected columns in the exact order defined by the output schema:

INSPECTION DATE, CAMIS, DBA, BORO, BUILDING, STREET, ZIPCODE, PHONE, CUISINE DESCRIPTION, ACTION, VIOLATION CODE, VIOLATION DESCRIPTION, CRITICAL FLAG, SCORE, GRADE, INSPECTION TYPE

This ensures that the output CSV can be directly imported into tools like Hive or Trino.

### Data Record Writing

The `reduce()` method receives grouped values from the mapper and writes them directly to the output.

- The reducer receives all mapper outputs grouped as an `Iterable` of records for each key. Since the mapper uses an empty key, each group effectively contains a sequence of fully formatted CSV rows. The reducer simply iterates through this iterable and writes each record directly to the final output without performing any additional transformations.

- **Output Format:**

  - **Key:** `NullWritable`, this ensures no key field appears in the final CSV output.
  - **Value:** The cleaned and processed CSV row generated by the mapper.

**Output location**

The output is stored in a single file because `setNumReduceTasks(1)` in my driver code collects all cleaned records into one reducer, producing a consolidated CSV output (`part-r-00000`) which is easy to load and analyse.

# Running MapReduce code

In this part now I am actually running the MapReduce code on the Dataproc Hadoop.

## Uploading the data to Hadoop

Dataproc documentation provides an easy way to upload datasets to GCP ingest bucket and then directly copy them into the Dataproc Hadoop Distributed File System (HDFS). After uploading the dataset to the GCS bucket, I used Hadoop's distcp tool to transfer the files into my HDFS directory.



Figure 4: Step to upload dataset on GCP

To move the dataset from GCS into HDFS, I used the following command:

```
hadoop distcp gs://nyu-dataproc-hdfs-ingest/nyc_inspection_group21 /user/an4758_nyu_edu
```

Figure 5: Downloading data from GCP to Hadoop

To organize and download the dataset, I listed my HDFS contents, moved the CSV file from its folder into my HDFS root directory, and then copied it to my local filesystem as a backup copy.



Figure 6: Renaming file and storing a copy on my local Dataproc FS

**Commands:**

```
hadoop fs -ls
hadoop fs -mv nyc_inspection_group21/NYC_Restaurent_Inscp_Data.csv /user/an4758_nyu_edu/
hadoop fs -ls
hadoop fs -get /user/an4758_nyu_edu/NYC_Restaurent_Inscp_Data.csv .
du -h NYC_Restaurent_Inscp_Data.csv
```

## Compiling the code

There are two ways to compile code, I have included a build.sh file and also directly used javac to execute the commands and create a JAR file.



Figure 7: Compiling MapReduce code

```
# Download the required dependency
wget https://repo1.maven.org/maven2/org/apache/commons/commons-csv/1.10.0/commons-csv-1.10.0.jar

# Compile files with Hadoop classpath and CSV JAR
javac -classpath ".:commons-csv-1.10.0.jar:`hadoop classpath`" *.java

# Copy compiled class and unpack third party dependency.
mkdir -p jar_test_dir
cp *.class jar_test_dir/
cd jar_test_dir
jar xf ../commons-csv-1.10.0.jar

cd ..
# Create a complete JAR.
jar cf nyc-inspection.jar -C jar_test_dir .

rm -rf jar_test_dir
rm -rf nyc-inspection.jar

# To BUILD using build.sh
chmod +x build.sh
./build.sh
```

The above commands will build the code by compiling Java source files and will pack the dependencies by generating a final `nyc-inspection.jar`, that I will use to run the MapReduce job.

8

## Running MapReduce job

Now that I have a Jar file, I will run the code on Dataproc and get the required results.



Figure 8: Running MapReduce job for NYC inspection data.



Figure 9: Job counters and output file.

My MapReduce job ran successfully, and I was able to see all my custom counters printed in the logs (Total Records, Borough-level information, Invalid ZIPCODES and much more). The reducer output is in the `part1_output` directory in HDFS, containing the result file.

## Output

In this step, I copied the processed output from HDFS into the local Dataproc FS and evaluated the rows to verify the MapReduce job output. As seen in the screenshot, I was able to clearly view the CSV header followed by properly formatted rows. This confirms that the reducer output is correct. After validating the dataset, I compressed all required source files, scripts, dependency JARs, and outputs into a `.tar.gz` format to prepare for my part1 submission.



Figure 10: Checking output, copying file, analysing file and storing it as csv.



Figure 11: Preparing the files to download

Finally, I downloaded the file opened the processed output CSV to verify if it is correctly pro-

cessed. As shown below, the file contains the correct header followed by the correct row entries generated by the reducer. I also can see that my Grade and Score values are properly processed.



| INSPECTION DATE | CAMIS | DBA | BORO | BUILDING | STREET | ZIPCODE | PHONE | CUISINE DESC | ACTION | VIOLATION CODE | VIOLATION DESCRIPTION | CRITICAL FLAG | SCORE | GRADE | INSPECTION TYPE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 03/14/2025 | 50156692 | MEMO'S RESTAURAN | Queens | 90-21 | 31 AVENUE | 11369 | 3478845775 | Spanish | Violations were cited in the fc | 28-06 | Contract with a pest management profes | Not Critical | 13 | A | Pre-permit (Operational) / Re-inspection |
| 02/21/2025 | 50142907 | ROSTICCERIA EVELIN | Brooklyn | 455 | MYRTLE AVENUE | 11205 | 6465529587 | Italian | Violations were cited in the fc | 10F | Non-food contact surface or equipment | Not Critical | 13 | A | Pre-permit (Operational) / Initial Inspection |
| 06/30/2025 | 50140743 | RED KUP | Manhattan | 701 | SAINT NICHOLAS AVE | 10031 | 9293019405 | Coffee/Tea | Violations were cited in the fc | 06D | Food contact surface not properly washe | Critical | 67 | C | Cycle Inspection / Initial Inspection |
| 07/15/2025 | 50110670 | ANYTIME BAR & BILLI | Manhattan | 112 | WEST 30 STREET | 10001 | 6466372967 | American | Violations were cited in the fc | 06D | Food contact surface not properly washe | Critical | 33 | C | Cycle Inspection / Initial Inspection |
| 01/16/2025 | 50105561 | DAVIDOVICH BAKERY | Manhattan | 79 | CLINTON STREET | 10002 | 5168280218 | Bakery Products | Violations were cited in the fc | 08A | Establishment is not free of harborage or | Not Critical | 29 | C | Cycle Inspection / Initial Inspection |
| 09/05/2024 | 50150570 | PANERA BREAD #6366 | Bronx | 46 | WEST 225 STREET | 10463 | 3473357139 | American | Violations were cited in the fc | 10F | Non-food contact surface or equipment | Not Critical | 25 | B | Pre-permit (Operational) / Initial Inspection |
| 04/27/2022 | 50117186 | XIN LA GONG FU | Queens | 35-48 | UNION STREET | 11354 | 9173316661 | Korean | Violations were cited in the fc | 05D | Hand washing facility not provided in or | Critical | 22 | B | Pre-permit (Operational) / Initial Inspection |
| 02/17/2023 | 50131655 | LIBERTY BAGELS | Manhattan | 32 | BROADWAY | 10004 | 5164583408 | Bagels/Pretzels | Violations were cited in the fc | 10A | Toilet facility not maintained or provided | Not Critical | 21 | B | Pre-permit (Non-operational) / Initial Inspection |
| 05/29/2024 | 50067257 | SAJHOMA RESTAURAI | Brooklyn | 408 | NEW LOTS AVENUE | 11207 | 3476737555 | Spanish | Violations were cited in the fc | 02H | After cooking or removal from hot holding | Critical | 19 | B | Cycle Inspection / Initial Inspection |
| 03/05/2025 | 50130900 | PELICANA CHICKEN | Queens | 47-08 | GREENPOINT AVENU | 11104 | 9293319793 | Chicken | Violations were cited in the fc | 02G | Cold TCS food item held above 41 °F; sr | Critical | 30 | C | Cycle Inspection / Re-inspection |
| 01/14/2022 | 40400209 | WHEELER'S | Brooklyn | 1705 | SHEEPSHEAD BAY R( | 11235 | 7186469320 | American | Violations were cited in the fc | 04H | Raw, cooked or prepared food is adultera | Critical | 19 | B | Cycle Inspection / Initial Inspection |
| 04/16/2025 | 40393093 | RINCON SALVADOREN | Queens | 92-15 | 149 STREET | 11435 | 5167325528 | Latin American | Violations were cited in the fc | 04K | Evidence of rats or live rats in establish | Critical | 43 | C | Cycle Inspection / Initial Inspection |
| 10/16/2023 | 50041617 | ITTADI GARDEN & GRI | Queens | 73-07 | 37 ROAD | 11372 | 3478662923 | Bangladeshi | Violations were cited in the fc | 08A | Establishment is not free of harborage or | Not Critical | 18 | B | Cycle Inspection / Re-inspection |
| 02/15/2024 | 50072519 | MIKE'S DINER | Brooklyn | 1454 | 86 STREET | 11228 | 3475679223 | Greek | Violations were cited in the fc | 08A | Establishment is not free of harborage or | Not Critical | 25 | B | Cycle Inspection / Re-inspection |
| 02/07/2025 | 50144653 | CAFE ON 7TH | Brooklyn | 493 | 7 AVENUE | 11215 | 3472357158 | American | Violations were cited in the fc | 04L | Evidence of mice or live mice in establish | Critical | 48 | C | Cycle Inspection / Initial Inspection |
| 10/25/2022 | 50104889 | DON PEPE TORTAS Y | Brooklyn | 3908 | 5 AVENUE | 11232 | 7184353326 | Mexican | Violations were cited in the fc | 02G | Cold TCS food item held above 41 °F; sr | Critical | 59 | C | Cycle Inspection / Initial Inspection |
| 09/03/2025 | 50139765 | BAOZI | Brooklyn | 5405 | 8 AVENUE | 11220 | 6464343333 | Chinese | Violations were cited in the fc | 05D | No hand washing facility in or adjacent tc | Critical | 43 | C | Cycle Inspection / Initial Inspection |
| 04/24/2025 | 50165223 | EATON CAFE | Queens | 89-08 | QUEENS BOULEVARI | 11373 | 3475070276 | Japanese | Violations were cited in the fc | 02G | Cold TCS food item held above 41 °F; sr | Critical | 42 | C | Pre-permit (Operational) / Initial Inspection |
| 06/12/2024 | 50118305 | TAJMAHAL RESTAURA | Brooklyn | 473 | MCDONALD AVENUE | 11218 | 3472401071 | Bangladeshi | Establishment Closed by DO | 04K | Evidence of rats or live rats in establishm | Critical | 93 | C | Cycle Inspection / Initial Inspection |
| 07/17/2023 | 50116612 | ALIDORO | Manhattan | 383 | WEST 31 STREET | 10001 | 6466882924 | Italian | Violations were cited in the fc | 02H | Filth flies or food/refuse/sewage associat | Critical | 23 | B | Pre-permit (Operational) / Initial Inspection |
| 06/20/2024 | 50121127 | BEAR DONUT | Manhattan | 40 | WEST 31 STREET | 10001 | 2013148342 | Donuts | Violations were cited in the fc | 08A | Establishment is not free of harborage or | Not Critical | 15 | B | Cycle Inspection / Initial Inspection |
| 08/21/2025 | 50127144 | IHOP | Brooklyn | 2951 | AVENUE U | 11229 | 7187192630 | Pancakes/Waffle | Violations were cited in the fc | 06F | Wiping cloths not stored clean and dry, c | Critical | 20 | B | Cycle Inspection / Initial Inspection |
| 05/29/2024 | 50129401 | GARDEN BAR & GRILL | Bronx | 3392 | EAST TREMONT AVE | 10461 | 7186640566 | Latin American | Violations were cited in the fc | 04L | Evidence of mice or live mice in establish | Critical | 34 | C | Cycle Inspection / Initial Inspection |
| 03/31/2025 | 50001870 | THE DUMPLING COVE | Bronx | 1530 | EAST 222 STREET | 10469 | 7186533143 | Caribbean | Violations were cited in the fc | 02B | Hot TCS food item not held at or above 1 | Critical | 35 | C | Cycle Inspection / Initial Inspection |
| 05/20/2024 | 50112033 | WANPO TEA SHOP | Manhattan | 37 | EAST 8 STREET | 10003 | 2129958349 | Coffee/Tea | Violations were cited in the fc | 02B | Hot TCS food item not held at or above 1 | Critical | 12 | A | Cycle Inspection / Initial Inspection |
| 11/13/2023 | 41304936 | DRAM SHOP | Brooklyn | 339 | 9 STREET | 11215 | 7187881444 | American | Violations were cited in the fc | 04H | Raw, cooked or prepared food is adultera | Critical | 42 | C | Cycle Inspection / Initial Inspection |
| 01/06/2025 | 50123802 | GOLDEN GATE EXPRE | Bronx | 300 | WEST 231 STREET | 10463 | 7188840077 | Chinese | Violations were cited in the fc | 10F | Non-food contact surface or equipment | Not Critical | 22 | B | Cycle Inspection / Initial Inspection |
| 04/30/2025 | 50066563 | RONI LIRA BROTHERS | Queens | 44-44 | COLLEGE POINT BOL | 11355 | 3476100497 | American | Violations were cited in the fc | 04L | Evidence of mice or live mice in establish | Critical | 27 | B | Cycle Inspection / Re-inspection |
| 10/21/2025 | 40679229 | AMARANTH | Manhattan | 21 | EAST 62 STREET | 10065 | 2129806700 | Mediterranean | Violations were cited in the fc | 10G | Dishwashing and ware washing: Cleaning | Not Critical | 13 | A | Cycle Inspection / Re-inspection |
| 03/16/2023 | 41563707 | PIZZERIA GIOVE | Staten Island | 278 | NEW DORP LANE | 10306 | 3472860635 | Italian | Violations were cited in the fc | 08C | Pesticide not properly labeled or used by | Not Critical | 9 | A | Cycle Inspection / Initial Inspection |
| 01/19/2022 | 41710752 | PALACE CAFE | Brooklyn | 2603 | NOSTRAND AVENUE | 11210 | 7183389525 | Jewish/Kosher | Violations were cited in the fc | 02G | Cold food item held above 41° F (smoked | Critical | 12 | A | Cycle Inspection / Re-inspection |
| 08/03/2023 | 50114993 | SIP SAK | Manhattan | 928 | 2 AVENUE | 10022 | 2125831900 | Turkish | Violations were cited in the fc | 02H | After cooking or removal from hot holding | Critical | 55 | C | Cycle Inspection / Initial Inspection |
| 07/18/2025 | 50064557 | VENIERO'S BAKERY | Manhattan | 340 | EAST 11 STREET | 10003 | 2126747070 | Bakery Products | Violations were cited in the fc | 10B | Anti-siphonage or back-flow prevention | Not Critical | 12 | A | Cycle Inspection / Initial Inspection |
| 10/13/2021 | 50107445 | HALAL BROS GRILL | Queens | 218-74 | HEMPSTEAD AVENUE | 11429 | 3479930857 | Chicken | Violations were cited in the fc | 06C | Food not protected from potential source | Critical | 10 | A | Pre-permit (Operational) / Initial Inspection |
| 06/26/2023 | 50101173 | KIKU SUSHI | Brooklyn | 453 | 7 AVENUE | 11215 | 7183691155 | Japanese | Violations were cited in the fc | 04M | Live roaches in facility's food or non-food | Critical | 28 | C | Cycle Inspection / Initial Inspection |
| 03/10/2023 | 50129871 | SUSHI D | Brooklyn | 207 | DEKALB AVENUE | 11205 | 7188580058 | Seafood | Violations were cited in the fc | 04E | Toxic chemical or pesticide improperly st | Critical | 28 | C | Pre-permit (Operational) / Initial Inspection |
| 03/13/2023 | 50133359 | AMMI | Manhattan | 25 | 11 AVENUE | 10011 | 2016967222 | Bangladeshi | Violations were cited in the fc | 06C | Food, supplies, or equipment not protect | Critical | 95 | C | Pre-permit (Non-operational) / Initial Inspection |
| 12/06/2024 | 50108155 | ABY'S BAR | Brooklyn | 1541 | MYRTLE AVENUE | 11237 | 3479442568 | Spanish | Violations were cited in the fc | 04A | Food Protection Certificate (FPC) not hel | Critical | 37 | N | Cycle Inspection / Initial Inspection |
| 12/27/2022 | 50056438 | KABAYAN FILIPINO RE | Queens | 69-12 | ROOSEVELT AVENUE | 11377 | 7182054010 | Filipino | Violations were cited in the fc | 02B | Hot TCS food item not held at or above 1 | Critical | 27 | B | Cycle Inspection / Re-inspection |
| 05/25/2023 | 50040547 | KITCHEN GRILL | Brooklyn | 914A | FULTON STREET | 11238 | 7187897800 | Indian | Violations were cited in the fc | 08A | Establishment is not free of harborage or | Not Critical | 16 | B | Cycle Inspection / Initial Inspection |
| 03/11/2025 | 50113746 | PANINO | Brooklyn | 5401 | 13 AVENUE | 11219 | 9292899013 | Jewish/Kosher | Establishment Closed by DO | 05D | No hand washing facility in or adjacent tc | Critical | 102 | N | Pre-permit (Operational) / Re-inspection |
| 12/17/2024 | 50063458 | RUMA'S KITCHEN | Queens | 37-01 | 61 STREET | 11377 | 7188999100 | Bangladeshi | Violations were cited in the fc | 08A | Establishment is not free of harborage or | Not Critical | 53 | C | Cycle Inspection / Re-inspection |
| 07/23/2022 | 50118545 | LA FLOR DEL PARAIS( | Queens | 80-29 | JAMAICA AVENUE | 11421 | 9292756890 | Spanish | Violations were cited in the fc | 02G | Cold TCS food item held above 41 °F; sr | Critical | 12 | A | Pre-permit (Operational) / Initial Inspection |
| 03/14/2022 | 50074586 | MIKE JR'S RICHMOND | Staten Island | 3954 | RICHMOND AVENUE | 10312 | 7183172331 | American | Violations were cited in the fc | 10B | Plumbing not properly installed or mainta | Not Critical | 28 | C | Cycle Inspection / Initial Inspection |

Figure 12: Finished processed output_data.csv file

# Conclusion

In this part of the project, I used MapReduce to process the NYC restaurant inspection dataset. I implemented validation, filtering, and column-level transformations in mapper and produced a processed and well-formatted CSV output using the reducer. I also collected specific information using specific counters. I also manually opened the processed csv file and checked it. This completes the data ingestion phase. The cleaned dataset is ready for further analysis using Hive, Trino and Tableau in the next part of the project.