

# IST 769 Final Project

Sana Khan

IST 769 Spring 2023

## Summary

The purpose of this project was to demonstrate skills in database management by using Spark to store time-series data in a Cassandra Keyspace and query it via Drill. Two different data sources were used to demonstrate proficiency in pulling data from an API. The first data source was NYC 311 complaints, which saved the data in a JSON. The second source was Weather data for Syracuse New York and the data was saved as a CSV. The steps and code are provided below.

## Steps:

1. Create a Keyspace in Cassandra
2. Pull data from the API and store it as a JSON/CSV
3. Read the JSON/CSV and store the data in a variable that's transferred to a Data Frame.
4. Create a table in Jupyter that is transcribed to a Cassandra Keyspace using Spark
5. Write the data in the data frame to the table created in step 4 in the keyspace
6. Enable the Cassandra plugin in Drill
7. Query the newly created tables in Drill

## NYC 311 complaints

Create a Cassandra Keyspace called `project_311`

```
1 CREATE KEYSPACE IF NOT EXISTS project311
2 WITH replication = {'class': 'SimpleStrategy', 'replication_factor': 1};
```

Pull the data from the API and save it as JSON/CSV

```
1 !pip install pandas
2 !pip install sodapy
3
4 import pandas as pd
5 from sodapy import Socrata
6
7 data_url='data.cityofnewyork.us'
8 data_set='erm2-nwe9'
9 app_token='T2Jeqm4J7zvAiJruefHcYvb7y'
10 selectc = "unique_key, created_date, agency, complaint_type, descriptor, location_type, intersection_street_1, i
11 client = Socrata(data_url,app_token)
12 client.timeout = 360
13
```

```

14 results = client.get(data_set,select=selectc, limit=2000)
15
16 df = pd.DataFrame.from_records(results)
17
18 df.to_json("my_311_data.json", orient="records")
19 df.to_csv("my_311_data.csv")

```

Read the JSON file and save it as a variable

```

1 file_path = 'file:///home/jovyan/datasets/my_311_data.json'
2 df311= spark.read.json(file_path)
3 df311.printSchema()
4
5 df311.count()
6 df311.select('complaint_type').distinct().count()
7

```

```
df311.count()
```

```
2000
```

```
df311.select('complaint_type').distinct().count()
```

```
81
```

Save the table in Cassandra

```

1 create_table_sql = '''
2 CREATE TABLE project311.data
3 (
4
5     address_type text,
6     agency text,
7     borough text,
8     city text,
9     community_board text,
10    complaint_type text,
11    created_date timestamp,
12    descriptor text,
13    incident_address text,
14    incident_zip text,
15    intersection_street_1 text,
16    intersection_street_2 text,
17    latitude float,
18    location_type text,
19    longitude float,
20    open_data_channel_type text,
21    resolution_action_updated_date timestamp,
22    status text,
23    unique_key text,
24    PRIMARY KEY (unique_key)

```

```

25
26 );
27 '''
28 from cassandra.cluster import Cluster
29 with Cluster([cassandra_host]) as cluster:
30     session = cluster.connect()
31 session.execute(create_table_sql)
32 print(create_table_sql)
33

```

```

CREATE TABLE project311.data
(
    address_type text,
    agency text,
    borough text,
    city text,
    community_board text,
    complaint_type text,
    created_date timestamp,
    descriptor text,
    incident_address text,
    incident_zip text,
    intersection_street_1 text,
    intersection_street_2 text,
    latitude float,
    location_type text,
    longitude float,
    open_data_channel_type text,
    resolution_action_updated_date timestamp,
    status text,
    unique_key text,
    PRIMARY KEY (unique_key)
);

```

Save the data from the JSON into the table `data`

```

1 df2 = df311.toDF(
2
3 "address_type",
4 "agency",
5 "borough",
6 "city",
7 "community_board",
8 "complaint_type",
9 "created_date",
10 "descriptor",
11 "incident_address",
12 "incident_zip",
13 "intersection_street_1",

```

```

14 "intersection_street_2",
15 "latitude",
16 "location_type",
17 "longitude",
18 "open_data_channel_type",
19 "resolution_action_updated_date",
20 "status",
21 "unique_key")
22
23 df2.write.format("org.apache.spark.sql.cassandra")\
24     .mode("Append")\
25     .option("table", "data")\
26     .option("keyspace", "project311")\
27     .save()
28

```

View data in Cassandra to confirm it has been added

```

1 select * from project311.data311;

```

```

cqlsh> select * from project311.data311;

```

unique_key	address_type	agency	borough	city	community_board	complaint_type	created_date	descriptor	latitude	location_type	longitude	open_data_channel_type
578463118	ADDRESS	NYPD	BRONX	10469	BRONX	11 BRONX	2023-06-09 22:59:38.000000+0000	Illegal Parking	40.86314	Street/Sidewalk	-73.83707	MOBILE
ked Hydrant	2515 LODOVICK AVENUE				MACE AVENUE	ALLERTON AVENUE						
57846412	ADDRESS	NYPD	MANHATTAN	10016	NEW YORK	05 MANHATTAN	2023-06-09 21:27:54.000000+0000	Noise - Commercial	40.74479	Club/Bar/Restaurant	-73.98182	Loud
Music/Party	118 EAST 31 STREET				PARK AVENUE SOUTH	LEXINGTON AVENUE						ONLINE
2023-06-09 21:06:56.000000+0000												
57847041	ADDRESS	NYPD	QUEENS	11373	ELMHURST	04 QUEENS	2023-06-09 21:59:28.000000+0000	Noise - Residential	40.74204	Residential Building/House	-73.88493	Loud
Music/Party	42-72 80 STREET				WOODSIDE AVENUE	45 AVENUE						ONLINE
2023-06-09 22:32:26.000000+0000												
57847389	ADDRESS	NYPD	BRONX	10462	BRONX	11 BRONX	2023-06-09 21:29:18.000000+0000	Illegal Parking	40.84683	Street/Sidewalk	-73.85513	Commercial Overn
ght Parking	1000 VAN NEST AVENUE				COLDEN AVENUE	PAULDING AVENUE						ONLINE
2023-06-09 22:28:01.000000+0000												

Query the data using Spark

```

1 df3 =spark.read.format("org.apache.spark.sql.cassandra")\
2
3     .option("table", "data")\
4
5     .option("keyspace", "project311")\
6
7     .load()
8
9 df3.count()

```

```

5]: 2000

6]: df3.printSchema()

root
|-- unique_key: string (nullable = false)
|-- address_type: string (nullable = true)
|-- agency: string (nullable = true)
|-- borough: string (nullable = true)
|-- city: string (nullable = true)
|-- community_board: string (nullable = true)
|-- complaint_type: string (nullable = true)
|-- created_date: timestamp (nullable = true)
|-- descriptor: string (nullable = true)
|-- incident_address: string (nullable = true)
|-- incident_zip: string (nullable = true)
|-- intersection_street_1: string (nullable = true)
|-- intersection_street_2: string (nullable = true)
|-- latitude: float (nullable = true)
|-- location_type: string (nullable = true)
|-- longitude: float (nullable = true)
|-- open_data_channel_type: string (nullable = true)
|-- resolution_action_updated_date: timestamp (nullable = true)
|-- status: string (nullable = true)

```

Query the data using Drill

```
1 SELECT * from Cassandra.project311.`data`
```

Show 10 entries

Search:  Show / hide columns

unique_key	address_type	agency	borough	city	community_board	complaint_type	created_date	descriptor
57846318	ADDRESS	NYPD	BRONX	BRONX	11 BRONX	Illegal Parking	1686351578000	Blocked Hydrant
57846412	ADDRESS	NYPD	MANHATTAN	NEW YORK	05 MANHATTAN	Noise - Commercial	1686346074000	Loud Music/Party
57847041	ADDRESS	NYPD	QUEENS	ELMHURST	04 QUEENS	Noise - Residential	1686347968000	Loud Music/Party
57847389	ADDRESS	NYPD	BRONX	BRONX	11 BRONX	Illegal Parking	1686346150000	Commercial Overnight Parking
57850517	ADDRESS	NYPD	QUEENS	BAYSIDE	11 QUEENS	Illegal Parking	1686346845000	Blocked Crosswalk
57850599	ADDRESS	NYPD	BROOKLYN	BROOKLYN	05 BROOKLYN	Blocked Driveway	1686359095000	Partial Access
57852550	ADDRESS	NYPD	MANHATTAN	NEW YORK	12 MANHATTAN	Noise - Street/Sidewalk	1686341579000	Loud Music/Party
57851442	ADDRESS	NYPD	BROOKLYN	BROOKLYN	04 BROOKLYN	Noise - Residential	1686358316000	Loud Music/Party
57844112	ADDRESS	NYPD	BRONX	BRONX	06 BRONX	Noise - Residential	1686360930000	Loud Music/Party
57846337	ADDRESS	NYPD	BROOKLYN	BROOKLYN	02 BROOKLYN	Illegal Parking	1686354300000	Blocked Hydrant

```

1 SELECT count(unique_key) as Num_Complaints,
2 borough
3 from Cassandra.project311.`data`
4 group by 2
5 order by 1 DESC
6

```

Num_Complaints	borough
575	BROOKLYN
558	QUEENS
413	BRONX
402	MANHATTAN
51	STATEN ISLAND
1	Unspecified

```

1 SELECT count(unique_key) as Num_Complaints,
2 complaint_type,
3 borough
4 from Cassandra.project311.`data`

```

```

5 group by 2,3
6 order by 3
7

```

Show 10 entries Search:  Show / hide columns

unique_key	address_type	agency	borough	city	community_board	complaint_type	created_date	descriptor
57846318	ADDRESS	NYPD	BRONX	BRONX	11 BRONX	Illegal Parking	1686351578000	Blocked Hydrant
57846412	ADDRESS	NYPD	MANHATTAN	NEW YORK	05 MANHATTAN	Noise - Commercial	1686346074000	Loud Music/Party
57847041	ADDRESS	NYPD	QUEENS	ELMHURST	04 QUEENS	Noise - Residential	1686347968000	Loud Music/Party
57847389	ADDRESS	NYPD	BRONX	BRONX	11 BRONX	Illegal Parking	1686346150000	Commercial Overnight Parking
57850517	ADDRESS	NYPD	QUEENS	BAYSIDE	11 QUEENS	Illegal Parking	1686346845000	Blocked Crosswalk
57850599	ADDRESS	NYPD	BROOKLYN	BROOKLYN	05 BROOKLYN	Blocked Driveway	1686359095000	Partial Access
57852550	ADDRESS	NYPD	MANHATTAN	NEW YORK	12 MANHATTAN	Noise - Street/Sidewalk	1686341579000	Loud Music/Party
57851442	ADDRESS	NYPD	BROOKLYN	BROOKLYN	04 BROOKLYN	Noise - Residential	1686358316000	Loud Music/Party
57844112	ADDRESS	NYPD	BRONX	BRONX	06 BRONX	Noise - Residential	1686360930000	Loud Music/Party
57846337	ADDRESS	NYPD	BROOKLYN	BROOKLYN	02 BROOKLYN	Illegal Parking	1686354300000	Blocked Hydrant

Showing 1 to 10 of 200 entries

## Open Meteo Weather Data

Create a Cassandra Keyspace called `project_weather_data`

```

1 CREATE KEYSPACE IF NOT EXISTS project_weather_data
2 WITH REPLICATION = {
3   'class': 'SimpleStrategy',
4   'replication_factor': 1
5 };

```

Code to pull data from the API:

```

1 import requests
2 response = requests.get('https://archive-api.open-meteo.com/v1/archive?latitude=43.0481&longitude=76.1474&start_c
3 data = response.json()

```

Code to save data as CSV

```

1 from pyspark.sql.functions import col,explode
2 import pandas as pd
3 df = pd.json_normalize(data)
4 df.to_json("project_weather_data_meteo.json2", orient="records")
5 df.to_csv("project_weather_data_meteo.csv2", index=False, header=True)

```

Code to open the CSV and save it as a data frame:

```

1
2 df_weather= pd.read_csv(file_weather)

```

```

3 df_weather = df_weather.rename(columns={'Time': 'hour'})
4
5 df_spark = spark.createDataFrame(df_weather)
6

```

```
df_spark.printSchema()
```

```

root
|-- Date: string (nullable = true)
|-- hour: string (nullable = true)
|-- time: string (nullable = true)
|-- temperature_2m: double (nullable = true)
|-- dewpoint_2m: double (nullable = true)
|-- precipitation: double (nullable = true)
|-- rain: double (nullable = true)
|-- snowfall: double (nullable = true)
|-- weathercode: long (nullable = true)
|-- cloudcover: long (nullable = true)
|-- cloudcover_low: long (nullable = true)
|-- windspeed_10m: double (nullable = true)
|-- windspeed_100m: double (nullable = true)
|-- winddirection_10m: long (nullable = true)
|-- winddirection_100m: long (nullable = true)
|-- windgusts_10m: double (nullable = true)

```

Code to create a table in Jupyter notebook under `project_weather_data` keyspace in Cassandra and save the data to table called `weather_data_table`

```

1 create_table_sql = '''
2 CREATE TABLE IF NOT EXISTS project_weather_data.weather_data_table (
3     date date,
4     hour time,
5     time timestamp,
6     temperature_2m float,
7     dewpoint_2m float,
8     precipitation float,
9     rain float,
10    snowfall float,
11    weathercode float,
12    cloudcover float,
13    cloudcover_low float,
14    windspeed_10m float,
15    windspeed_100m float,
16    winddirection_10m float,
17    winddirection_100m float,
18    windgusts_10m float,
19    PRIMARY KEY (time)
20 );
21 '''

```

```

22 from cassandra.cluster import Cluster
23 with Cluster([cassandra_host]) as cluster:
24     session = cluster.connect()
25     session.execute(create_table_sql)
26 print(create_table_sql)
27

```

Code to write the data back to Cassandra

```

1  weather_2 = df_spark.toDF("date",
2  "hour",
3  "time",
4  "temperature_2m",
5  "dewpoint_2m",
6  "precipitation",
7  "rain",
8  "snowfall",
9  "weathercode",
10 "cloudcover",
11 "cloudcover_low",
12 "windspeed_10m",
13 "windspeed_100m",
14 "winddirection_10m",
15 "winddirection_100m",
16 "windgusts_10m");
17
18 weather_2.write.format("org.apache.spark.sql.cassandra")\
19     .mode("Append")\
20     .option("table", "weather_data_table")\
21     .option("keyspace", "project_weather_data")\
22     .save()

```

Code to query the table in cql

```

1 project_weather_data> select * from project_weather_data.weather_data_table;

```

Code to query in Spark

```

1 weather_3 =spark.read.format("org.apache.spark.sql.cassandra")\
2     .option("table", "weather_data_table")\
3     .option("keyspace", "project_weather_data")\
4     .load()
5
6 weather_3.count()

```

```

: weather_3 =spark.read.format("org.apache.spark.sql.cassandra")\
    .option("table", "weather_data_table")\
    .option("keyspace", "project_weather_data")\
    .load()

weather_3.count()

```

```

: 360

```



```

1 weather_3.createOrReplaceTempView("weather")
2 query = '''
3     SELECT time, temperature_2m, precipitation, rain, snowfall, weathercode FROM weather_data_table;
4     '''
5 spark.sql(query).explain()
6

```

```

15]: weather_3.createOrReplaceTempView("weather")
query = '''
SELECT time, temperature_2m, precipitation, rain, snowfall, weathercode FROM weather_data_table;
'''
spark.sql(query).explain()

== Physical Plan ==
*(1) Project [time#123, temperature_2m#132, precipitation#129, rain#130, snowfall#131, weathercode#133]
+- BatchScan[time#123, precipitation#129, rain#130, snowfall#131, temperature_2m#132, weathercode#133] Cassandra Scan: project_weather_data.weather_data_table
   - Cassandra Filters: []
   - Requested Columns: [time,precipitation,rain,snowfall,temperature_2m,weathercode]

```

```

1 query = '''
2     SELECT time, temperature_2m, precipitation, rain, snowfall, weathercode from weather_data_table WHERE temperature_2m > 60
3     '''
4 spark.sql(query).explain()
5

```

```

query = '''
SELECT time, temperature_2m, precipitation, rain, snowfall, weathercode from weather_data_table WHERE temperature_2m > 60
'''
spark.sql(query).explain()

== Physical Plan ==
*(1) Project [time#123, temperature_2m#132, precipitation#129, rain#130, snowfall#131, weathercode#133]
+- *(1) Filter (temperature_2m#132 > 60.0)
   +- BatchScan[time#123, precipitation#129, rain#130, snowfall#131, temperature_2m#132, weathercode#133] Cassandra Scan: project_weather_data.weather_data_table
      - Cassandra Filters: []
      - Requested Columns: [time,precipitation,rain,snowfall,temperature_2m,weathercode]

```

## Query in Drill

```

1 SELECT * from Cassandra.project_weather_data.`weather_data_table`

```

time	cloudcover	cloudcover_low	data	deeppoint_2m	hour	precipitation	rain	snowfall	temperature_2m	weathercode	winddirection_100m	winddirection_15m
684906000000	80.0	2.0	19501	30.2	2760000000000000	0.1	0.1	0.0	49.8	51.0	311.0	301.0
685408400000	1.0	1.0	19507	37.1	3600000000000000	0.0	0.0	0.0	42.6	0.0	141.0	145.0
684875600000	60.0	14.0	19500	38.6	7560000000000000	0.0	0.0	0.0	41.4	2.0	9.0	149.0
684918800000	85.0	15.0	19501	29.7	3240000000000000	0.0	0.0	0.0	51.8	3.0	315.0	288.0
685394000000	0.0	0.0	19506	31.5	7560000000000000	0.0	0.0	0.0	39.7	0.0	99.0	127.0
684881200000	70.0	23.0	19500	39.2	6120000000000000	0.0	0.0	0.0	44.0	2.0	207.0	159.0
684817640000	6.0	0.0	19502	33.9	3600000000000000	0.0	0.0	0.0	39.4	0.0	108.0	150.0
685473200000	20.0	0.0	19507	37.1	6840000000000000	0.0	0.0	0.0	44.4	0.0	81.0	113.0
684951200000	7.0	0.0	19501	30.5	6480000000000000	0.0	0.0	0.0	38.3	0.0	198.0	160.0
685044800000	5.0	0.0	19502	33.2	7200000000000000	0.0	0.0	0.0	38.2	0.0	202.0	166.0

```

1
2 SELECT * from cassandra.project_weather_data.`weather_data_table` where temperature_2m > '60'
3

```

time	cloudcover	cloudcover_low	data	dewpoint_2m	hour	precipitation	rain	snowfall	temperature_2m	weathercode	winddirection_100m	winddirection_1
1056137630000	43.0	0.0	19496	45.3	4600000000000000	0.0	0.0	0.0	66.0	1.0	270.0	270.0
10561377300000	28.0	0.0	19496	45.2	3960000000000000	0.0	0.0	0.0	68.7	1.0	5.0	360.0
10561378300000	32.0	1.0	19497	45.6	2520000000000000	0.4	0.4	0.0	64.3	51.0	300.0	310.0
10564940300000	36.0	1.0	19496	42.2	3960000000000000	0.0	0.0	0.0	60.3	1.0	217.0	214.0
10564942300000	71.0	2.0	19497	47.2	4680000000000000	0.3	0.3	0.0	61.4	51.0	200.0	193.0
10565264300000	18.0	7.0	19505	37.6	3240000000000000	0.1	0.1	0.0	60.2	51.0	300.0	301.0
10565436300000	0.0	0.0	19507	36.4	2760000000000000	0.0	0.0	0.0	60.8	0.0	6.0	3.0
10565734900000	38.0	4.0	19493	43.4	3960000000000000	0.1	0.1	0.0	63.7	51.0	240.0	236.0
10567185200000	43.0	3.0	19504	35.2	3960000000000000	0.3	0.3	0.0	63.3	51.0	250.0	260.0
10568653200000	28.0	1.0	19498	40.5	3600000000000000	0.0	0.0	0.0	62.0	1.0	282.0	286.0

Showing 1 to 10 of 45 entries

Previous12345Next

Conclusion

Cassandra is well suited for time series data given that it can handle writes with high reliability. Since the weather data has hourly data, Cassandra can handle the amount of data with ease. The 311 data also had a large volume of data, but the scope was limited to 2000 rows. The next iteration of this project would be to use the 311 data and show the complaints on a map using Elastic. I attempted to do this but the latitude/longitude was being stored in Elastic as a text variable even after transformations, which is incompatible with being able to be used on a map.