

## **P3: Implementation**

# **Boston Blue Bikes Analysis**

## **Project Group 15**

**Anjali Kabra**

**Jhalak Surve**

**Krishna Kapadia**

**Shubham Shah**

The objective of our database project is to analyze the blue bikes data to get some insights about the most used stations, most frequent routes, trip durations and the overall usage of the bikes. A few of the many questions we would like to answer through our analysis would be:

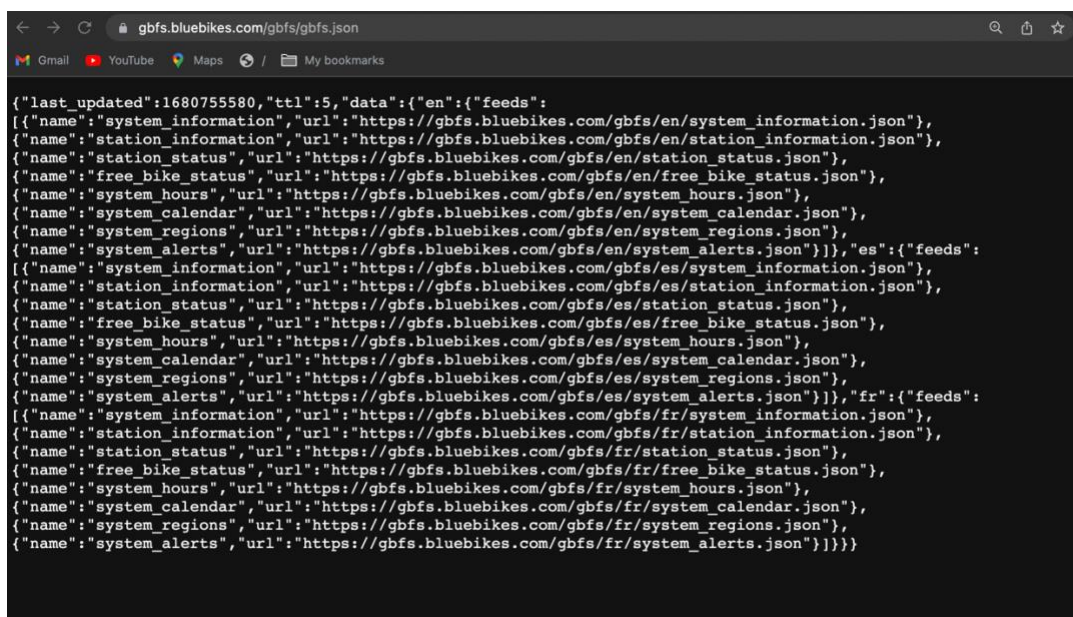
- What is the average duration for which a rider rents a bike?
- What is the most frequently accessed station?
- What are the rush hours?
- What are the busy blue bike days in the week?
- What is the ride count against the riders' age and gender?

DATA MODEL: ArangoDB Multi-Model (Graph + Document + Key value)

PLATFORM: Local

### Brief description of the implementation process:

- For the implementation, we have created a **data pipeline using python** to load the data into our database. (I have also submitted the jupyter notebook file)
- We are using the data from the Boston Blue Bikes official website.  
<https://www.bluebikes.com/system-data>
- Bluebikes publishes real-time system data in open General Bikeshare Feed Specification (GBFS) format. <https://gbfs.bluebikes.com/gbfs/gbfs.json>



The screenshot shows a web browser window with the address bar displaying 'gbfs.bluebikes.com/gbfs/gbfs.json'. The browser's address bar also shows search, share, and star icons. Below the address bar, there are several tabs: 'Gmail', 'YouTube', 'Maps', and 'My bookmarks'. The main content area of the browser displays a large block of JSON data, which is the GBFS feed specification. The JSON is structured as follows: it starts with a 'last\_updated' field (1680755580), a 'ttl' field (5), and a 'data' field. The 'data' field contains two objects: 'en' and 'fr'. Each of these objects contains a 'feeds' array. The 'en' array lists various system information feeds (system\_information, station\_information, station\_status, free\_bike\_status, system\_hours, system\_calendar, system\_regions, system\_alerts) with their respective URLs. The 'fr' array lists the same feeds for French. The JSON is formatted with indentation and line breaks for readability.

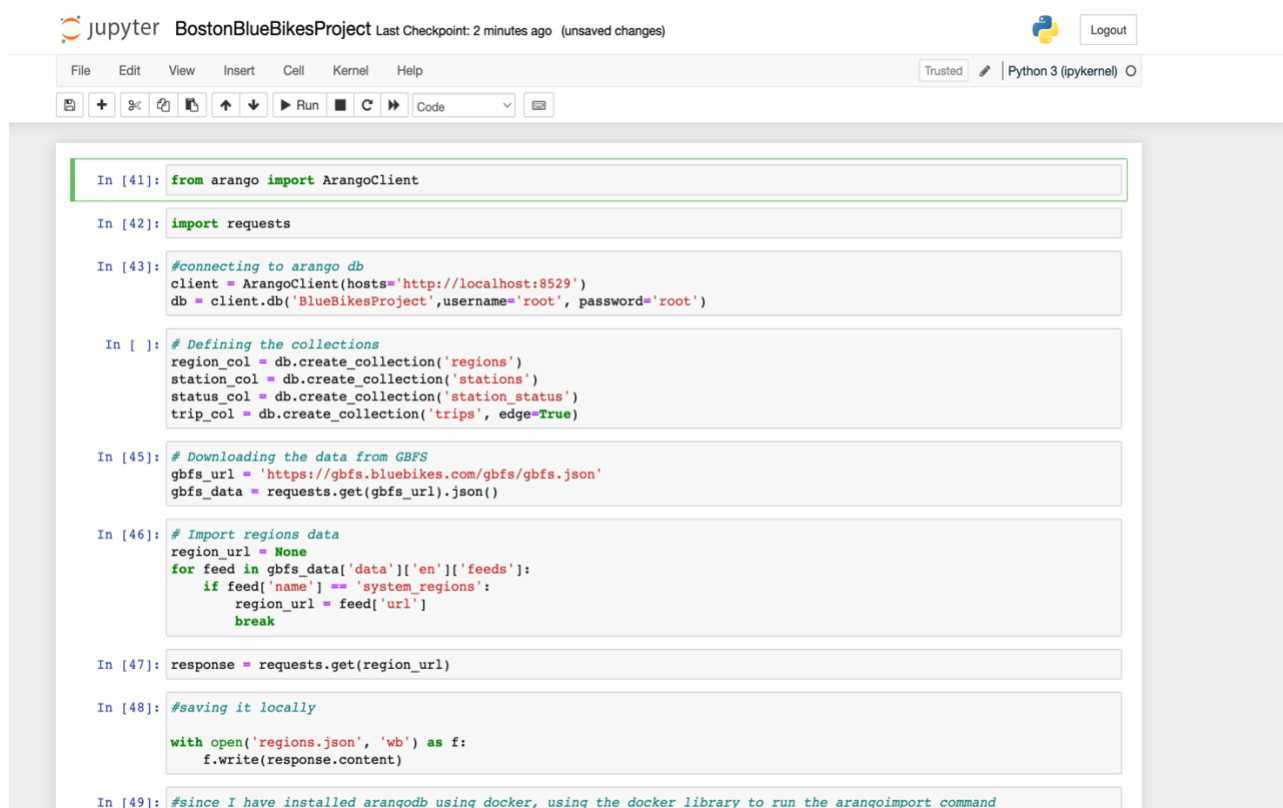
```
{
  "last_updated": 1680755580,
  "ttl": 5,
  "data": {
    "en": {
      "feeds": [
        {
          "name": "system_information",
          "url": "https://gbfs.bluebikes.com/gbfs/en/system_information.json"
        },
        {
          "name": "station_information",
          "url": "https://gbfs.bluebikes.com/gbfs/en/station_information.json"
        },
        {
          "name": "station_status",
          "url": "https://gbfs.bluebikes.com/gbfs/en/station_status.json"
        },
        {
          "name": "free_bike_status",
          "url": "https://gbfs.bluebikes.com/gbfs/en/free_bike_status.json"
        },
        {
          "name": "system_hours",
          "url": "https://gbfs.bluebikes.com/gbfs/en/system_hours.json"
        },
        {
          "name": "system_calendar",
          "url": "https://gbfs.bluebikes.com/gbfs/en/system_calendar.json"
        },
        {
          "name": "system_regions",
          "url": "https://gbfs.bluebikes.com/gbfs/en/system_regions.json"
        },
        {
          "name": "system_alerts",
          "url": "https://gbfs.bluebikes.com/gbfs/en/system_alerts.json"
        }
      ]
    },
    "fr": {
      "feeds": [
        {
          "name": "system_information",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/system_information.json"
        },
        {
          "name": "station_information",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/station_information.json"
        },
        {
          "name": "station_status",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/station_status.json"
        },
        {
          "name": "free_bike_status",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/free_bike_status.json"
        },
        {
          "name": "system_hours",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/system_hours.json"
        },
        {
          "name": "system_calendar",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/system_calendar.json"
        },
        {
          "name": "system_regions",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/system_regions.json"
        },
        {
          "name": "system_alerts",
          "url": "https://gbfs.bluebikes.com/gbfs/fr/system_alerts.json"
        }
      ]
    }
  }
}
```

- We are using this GBFS feed to fetch data related to the regions, stations and station status. For this purpose, we've used python.

Firstly, we have used the **python library 'arango'** to connect to the arangodb using python.

Then, we have defined four collections – **regions, stations, station\_status and trips(edge)**

Then, we are fetching the data using **requests.get()** method and saving the json file locally.



The image shows a Jupyter Notebook interface for a project named 'BostonBlueBikesProject'. The notebook contains the following code cells:

```
In [41]: from arango import ArangoClient

In [42]: import requests

In [43]: #connecting to arango db
client = ArangoClient(hosts='http://localhost:8529')
db = client.db('BlueBikesProject', username='root', password='root')

In [ ]: # Defining the collections
region_col = db.create_collection('regions')
station_col = db.create_collection('stations')
status_col = db.create_collection('station_status')
trip_col = db.create_collection('trips', edge=True)

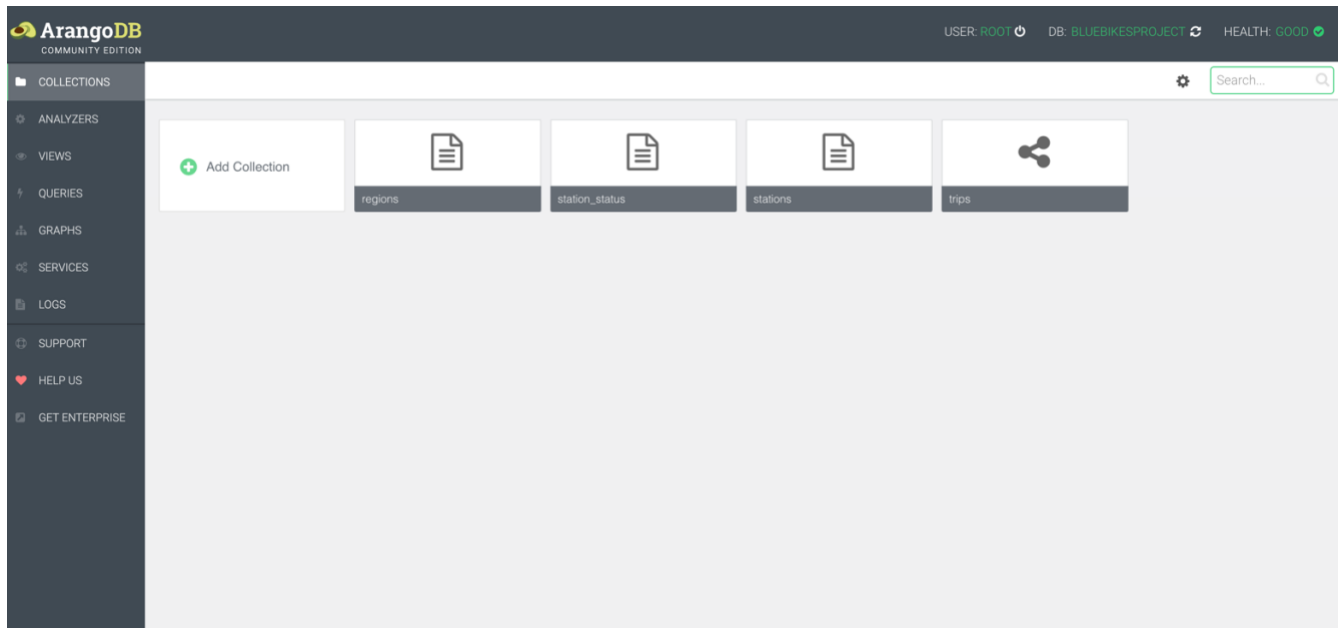
In [45]: # Downloading the data from GBFS
gbfs_url = 'https://gbfs.bluebikes.com/gbfs/gbfs.json'
gbfs_data = requests.get(gbfs_url).json()

In [46]: # Import regions data
region_url = None
for feed in gbfs_data['data']['en']['feeds']:
    if feed['name'] == 'system_regions':
        region_url = feed['url']
        break

In [47]: response = requests.get(region_url)

In [48]: #saving it locally
with open('regions.json', 'wb') as f:
    f.write(response.content)

In [49]: #since I have installed arangodb using docker, using the docker library to run the arangoimport command
```



- Then, we have used the **arangoinport** tool to import the json data into collections. We have written python code for the same.

jupyter BostonBlueBikesProject Last Checkpoint: 3 minutes ago (unsaved changes) Python 3 (ipykernel)

```

for feed in gis_data['data']:
    if feed['name'] == 'system_regions':
        region_url = feed['url']
        break

In [47]: response = requests.get(region_url)

In [48]: #saving it locally
with open('regions.json', 'wb') as f:
    f.write(response.content)

In [49]: #since I have installed arangodb using docker, using the docker library to run the arangoinport command
import docker
import os

In [50]: client = docker.from_env()

In [51]: container = client.containers.get('5aac08e9d503')

In [57]: #running the arangoinport code for importing the regions data
cmd = "arangoinport --file regions.json --type json --collection regions --server.endpoint tcp://{ip}:8529 --server.use
result = container.exec_run(cmd)
print(result.output.decode())

Connected to ArangoDB 'http+tcp://172.17.0.2:8529, version: 3.10.4, database: 'BlueBikesProject', username: 'root'
-----
database:          BlueBikesProject
collection:        regions
overwrite coll. prefix: no
create:            no
create database:   no
source filename:   regions.json
file type:         json
threads:           4
  
```

## Importing regions data:

```
jupyter BostonBlueBikesProject Last Checkpoint: 2 hours ago (unsaved changes) Python 3 (ipykernel)

File Edit View Insert Cell Kernel Help Trusted Python 3 (ipykernel)

In [57]: cmd = "arangosh --file regions.json --type json --collection regions --server.endpoint tcp://{ip}:8529 --server.username root
result = container.exec_run(cmd)
print(result.output.decode())

Connected to ArangoDB 'http+tcp://172.17.0.2:8529, version: 3.10.4, database: 'BlueBikesProject', username: 'root'
-----
database:           BlueBikesProject
collection:         regions
overwrite coll. prefix: no
create:             no
create database:    no
source filename:    regions.json
file type:          json
threads:            4
on duplicate:       error
connect timeout:    5
request timeout:    1200
-----
Starting JSON import...
2023-04-06T02:59:54Z [305] INFO [9ddf3] {general} processed 1.3 KB (3%) of input file

created:           20
warnings/errors:   0
updated/replaced:  0
ignored:           0
```

ArangoDB Collection: regions USER: root DB: BLUEBIKESPROJECT HEALTH: GOOD

Content Indexes Info Settings Computed Values Schema 10 results

Content	_key
{"_id":"regions/4","_key":"4","_rev":"_fzQ_mii--","name":"Brookline"}	4
{"_id":"regions/8","_key":"8","_rev":"_fzQ_mim--","name":"Cambridge"}	8
{"_id":"regions/9","_key":"9","_rev":"_fzQ_mim--","name":"Somerville"}	9
{"_id":"regions/10","_key":"10","_rev":"_fzQ_mim--A","name":"Boston"}	10
{"_id":"regions/33","_key":"33","_rev":"_fzQ_mim--B","name":"Motive"}	33
{"_id":"regions/39","_key":"39","_rev":"_fzQ_mim--C","name":"80"}	39
{"_id":"regions/89","_key":"89","_rev":"_fzQ_mim--D","name":"Hingham"}	89
{"_id":"regions/90","_key":"90","_rev":"_fzQ_miq--","name":"Quincy"}	90
{"_id":"regions/91","_key":"91","_rev":"_fzQ_miq--","name":"Hingham"}	91
{"_id":"regions/104","_key":"104","_rev":"_fzQ_miq--A","name":"Everett"}	104

20 doc(s) 1 2

## Importing stations data:

```
jupyter BostonBlueBikesProject Last Checkpoint: 5 hours ago (unsaved changes)
Python 3 (ipykernel)

In [59]: response = requests.get(station_url)

In [60]: #saving it locally
with open('stations.json', 'wb') as f:
    f.write(response.content)

In [61]: #running the arangoimport code for importing the stations data
cmd = "arangoimport --file stations.json --type json --collection stations --server.endpoint tcp://{ip}:8529 --server.u
result = container.exec_run(cmd)
print(result.output.decode())

Connected to ArangoDB 'http+tcp://172.17.0.2:8529, version: 3.10.4, database: 'BlueBikesProject', username: 'root'

-----
database:      BlueBikesProject
collection:    stations
overwrite coll. prefix: no
create:        no
create database: no
source filename: stations.json
file type:     json
threads:       4
on duplicate:  error
connect timeout: 5
request timeout: 1200
-----

Starting JSON import...
2023-04-06T03:11:31Z [317] INFO [9ddf3] {general} processed 288.8 KB (3%) of input file

created:      453
warnings/errors: 0
updated/replaced: 0
ignored:      0
```

ArangoDB Collection: stations

USER: root DB: BLUEBIKESPROJECT HEALTH: GOOD

Content Indexes Info Settings Computed Values Schema

10 results

Content	_key
{"_id": "stations/3", "_key": "3", "_rev": "_fzQKPge---", "capacity": 15, "eightd_has_key_dispenser": true, "eightd_station_s...	3
{"_id": "stations/4", "_key": "4", "_rev": "_fzQKPge--", "capacity": 19, "eightd_has_key_dispenser": false, "eightd_station_...	4
{"_id": "stations/5", "_key": "5", "_rev": "_fzQKPgi---", "capacity": 15, "eightd_has_key_dispenser": false, "eightd_station_...	5
{"_id": "stations/6", "_key": "6", "_rev": "_fzQKPgi--", "capacity": 15, "eightd_has_key_dispenser": false, "eightd_station_...	6
{"_id": "stations/7", "_key": "7", "_rev": "_fzQKPgi--A", "capacity": 15, "eightd_has_key_dispenser": false, "eightd_station_...	7
{"_id": "stations/8", "_key": "8", "_rev": "_fzQKPgi--B", "capacity": 19, "eightd_has_key_dispenser": false, "eightd_station_...	8
{"_id": "stations/9", "_key": "9", "_rev": "_fzQKPgi--C", "capacity": 15, "eightd_has_key_dispenser": false, "eightd_station_...	9
{"_id": "stations/10", "_key": "10", "_rev": "_fzQKPGm---", "capacity": 11, "eightd_has_key_dispenser": false, "eightd_statio...	10
{"_id": "stations/11", "_key": "11", "_rev": "_fzQKPGm--", "capacity": 15, "eightd_has_key_dispenser": false, "eightd_statio...	11
{"_id": "stations/12", "_key": "12", "_rev": "_fzQKPGm--A", "capacity": 26, "eightd_has_key_dispenser": true, "eightd_station_...	12

453 doc(s)

## Importing station\_status data:

```
jupyter BostonBlueBikesProject Last Checkpoint: 5 hours ago (unsaved changes) Python 3 (ipykernel)

File Edit View Insert Cell Kernel Help Trusted Python 3 (ipykernel)

stationstatus_url = feed['url']
break

In [67]: response = requests.get(stationstatus_url)

In [68]: #saving it locally
with open('station_status.json', 'wb') as f:
    f.write(response.content)

In [69]: #running the arangoimport code for importing the station_status data
cmd = "arangoimport --file station_status.json --type json --collection station_status --server.endpoint tcp://{ip}:8529"
result = container.exec_run(cmd)
print(result.output.decode())

Connected to ArangoDB 'http+tcp://172.17.0.2:8529, version: 3.10.4, database: 'BlueBikesProject', username: 'root'

-----
database:      BlueBikesProject
collection:    station_status
overwrite coll. prefix: no
create:        no
create database: no
source filename: station_status.json
file type:     json
threads:       4
on duplicate:  error
connect timeout: 5
request timeout: 1200
-----

Starting JSON import...
2023-04-06T03:18:33Z [329] INFO [9ddf3] {general} processed 199.5 KB (3%) of input file

created:      453
warnings/errors: 0
updated/replaced: 0
ignored:      0
```

ArangoDB Collection: station\_status USER: ROOT DB: BLUEBIKESPROJECT HEALTH: GOOD

Content Indexes Info Settings Computed Values Schema 10 results

Content	_key
{"_id": "station_status/728308", "_key": "728308", "_rev": "_fzQQrl6--", "eightd_has_available_keys": false, "is_installed..	728308
{"_id": "station_status/728309", "_key": "728309", "_rev": "_fzQQrm---", "eightd_has_available_keys": false, "is_installed..	728309
{"_id": "station_status/728310", "_key": "728310", "_rev": "_fzQQrm---", "eightd_has_available_keys": false, "is_installed..	728310
{"_id": "station_status/728311", "_key": "728311", "_rev": "_fzQQrm---A", "eightd_has_available_keys": false, "is_installed..	728311
{"_id": "station_status/728312", "_key": "728312", "_rev": "_fzQQrm---B", "eightd_has_available_keys": false, "is_installed..	728312
{"_id": "station_status/728313", "_key": "728313", "_rev": "_fzQQrm---C", "eightd_has_available_keys": false, "is_installed..	728313
{"_id": "station_status/728314", "_key": "728314", "_rev": "_fzQQrmC--", "eightd_has_available_keys": false, "is_installed..	728314
{"_id": "station_status/728315", "_key": "728315", "_rev": "_fzQQrmC--", "eightd_has_available_keys": false, "is_installed..	728315
{"_id": "station_status/728316", "_key": "728316", "_rev": "_fzQQrmC--A", "eightd_active_station_services": [{"id": "f8346a..	728316
{"_id": "station_status/728317", "_key": "728317", "_rev": "_fzQQrmC--B", "eightd_has_available_keys": false, "is_installed..	728317

453 doc(s)



- Now, we're not getting the trips data from the GBFS feed, so we have used **existing datasets for trips data** from the blue bikes official website, which is in **csv format**.
- We've **cleaned** this csv data using **python pandas library**.

Jupyter BostonBlueBikesProject Last Checkpoint: 5 hours ago (autosaved)

File Edit View Insert Cell Kernel Help Trusted Python 3 (ipykernel)

```

In [70]: #cleaning the trip data in the csv form
import pandas as pd

In [71]: df = pd.read_csv('tripdata022023.csv')

In [72]: #checking for missing values
df.isnull().sum()

Out[72]: tripduration      0
starttime      0
stoptime      0
start station id      0
start station name      0
start station latitude      0
start station longitude      0
end station id      0
end station name      0
end station latitude      0
end station longitude      0
bikeid      0
usertype      0
postal code      9696
dtype: int64

In [ ]:

```

Jupyter BostonBlueBikesProject Last Checkpoint: 5 hours ago (unsaved changes)

File Edit View Insert Cell Kernel Help Trusted Python 3 (ipykernel)

```

In [78]: #removing null values from postal code column
df['postal code'] = df['postal code'].fillna(0)

In [79]: #since 'trips' is an edge collection, renaming the start_station_id and end_station_id columns to _from and _to respectively
df = df.rename(columns={'start_station_id': '_from'})
df = df.rename(columns={'end_station_id': '_to'})

In [82]: #appending 'stations/' in front of _from and _to station ids so that they refer to the ids of the station collection
df['_from'] = df['_from'].astype(str)
df['_to'] = df['_to'].astype(str)
df['_from'] = df['_from'].apply(lambda x: 'stations/' + x)
df['_to'] = df['_to'].apply(lambda x: 'stations/' + x)

In [83]: df

Out[83]:

```

	tripduration	starttime	stoptime	_from	start station name	start station latitude	start station longitude	_to	end station name	end station latitude	end station longitude	bikeid	usertype
0	263	2023-02-01 00:00:50.0410	2023-02-01 00:05:13.7090	stations/330	30 Dane St	42.381001	-71.104025	stations/110	Harvard University Gund Hall at Quincy St / Ki...	42.376369	-71.114025	7334	Subs...
1	447	2023-02-01 00:02:50.1390	2023-02-01 00:10:17.4090	stations/413	Kennedy-Longfellow School 158 Spring St	42.369553	-71.085790	stations/386	Sennott Park Broadway at Norfolk Street	42.368605	-71.099302	3257	Subs...
2	302	2023-02-01 00:09:00.7270	2023-02-01 00:14:03.2560	stations/554	Forsyth St at Huntington Ave	42.339202	-71.090511	stations/27	Roxbury Crossing T Stop - Columbus Ave at	42.331184	-71.095171	7824	Subs...



- After cleaning the data, we have imported it using the **arangoimport** tool.

```

jupyter BostonBlueBikesProject Last Checkpoint: 5 hours ago (unsaved changes)
File Edit View Insert Cell Kernel Help Trusted Python 3 (pykernel)
In [85]: #running the arangoimport code for importing the trips data
cmd = "arangoimport --file tripdata022023_clean.csv --type csv --collection trips --server.endpoint tcp://(ip):8529 --s
result = container.exec_run(cmd)
print(result.output.decode())

-----
Connected to ArangoDB 'http+tcp://172.17.0.2:8529, version: 3.10.4, database: 'BlueBikesProject', username: 'root'
-----
database:      BlueBikesProject
collection:    trips
overwrite coll. prefix: no
create:        no
create database: no
source filename: tripdata022023_clean.csv
file type:     csv
quote:         "
separator:     ,
headers file:
threads:       4
on duplicate:  error
connect timeout: 5
request timeout: 1200
-----
Starting CSV import...
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 1.3 MB (3%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 2.3 MB (6%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 3.4 MB (9%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 4.4 MB (12%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 5.5 MB (15%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 6.5 MB (18%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 7.6 MB (21%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 8.6 MB (24%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 9.6 MB (27%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 10.7 MB (30%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 11.7 MB (33%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 12.8 MB (36%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 13.8 MB (39%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 14.9 MB (42%) of input file
2023-04-06T03:51:14Z [341] INFO [9ddf3] {general} processed 15.9 MB (45%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 17.0 MB (48%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 18.0 MB (51%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 19.1 MB (54%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 20.1 MB (57%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 21.2 MB (60%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 22.2 MB (63%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 23.3 MB (66%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 24.3 MB (69%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 25.4 MB (72%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 26.4 MB (75%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 27.5 MB (78%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 28.5 MB (81%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 29.6 MB (84%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 30.6 MB (87%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 31.7 MB (90%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 32.7 MB (93%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 33.8 MB (96%) of input file
2023-04-06T03:51:15Z [341] INFO [9ddf3] {general} processed 34.8 MB (99%) of input file

created:      152975
warnings/errors: 0
updated/replaced: 0
ignored:      0
lines read:   152977

```

ArangoDB Collection: trips

USER: root DB: BLUEBIKESPROJECT HEALTH: GOOD

Content Indexes Info Settings Computed Values Schema

10 results

Content	_key
{"_from": "stations/330", "_id": "trips/729602", "_key": "729602", "_rev": "_fzQumze---", "_to": "stations/110", "bikeid": "733...	729602
{"_from": "stations/413", "_id": "trips/729603", "_key": "729603", "_rev": "_fzQumzi---", "_to": "stations/386", "bikeid": "325...	729603
{"_from": "stations/554", "_id": "trips/729604", "_key": "729604", "_rev": "_fzQumzm---", "_to": "stations/27", "bikeid": "7824...	729604
{"_from": "stations/87", "_id": "trips/729605", "_key": "729605", "_rev": "_fzQumzq---", "_to": "stations/178", "bikeid": "8266...	729605
{"_from": "stations/55", "_id": "trips/729606", "_key": "729606", "_rev": "_fzQum0C---", "_to": "stations/32", "bikeid": "7431,...	729606
{"_from": "stations/583", "_id": "trips/729607", "_key": "729607", "_rev": "_fzQum0C---", "_to": "stations/386", "bikeid": "681...	729607
{"_from": "stations/131", "_id": "trips/729608", "_key": "729608", "_rev": "_fzQum0G---", "_to": "stations/125", "bikeid": "693...	729608
{"_from": "stations/362", "_id": "trips/729609", "_key": "729609", "_rev": "_fzQum0K---", "_to": "stations/362", "bikeid": "592...	729609
{"_from": "stations/24", "_id": "trips/729610", "_key": "729610", "_rev": "_fzQum0K---", "_to": "stations/58", "bikeid": "3085,...	729610
{"_from": "stations/358", "_id": "trips/729611", "_key": "729611", "_rev": "_fzQum00---", "_to": "stations/84", "bikeid": "7762...	729611

152,975 edge(s)