

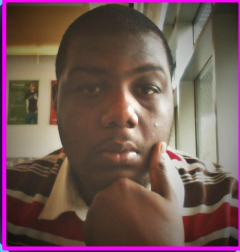
CS 488/588: Cloud & Cluster Data Management

Freeway Data in MongoDB

««« Portland State University, Fall 2020 »»»



The DJs



DOMINIQUE
MOORE



JANE
SEIGMAN



SANTIAGO
TOBON



TABLE OF CONTENTS



01

DATA

How we modeled
the freeway data in
MongoDB
&
ETL process

02

QUERIES

Query Demos
&
Query Plans

03

SELF-CRITIQUE

Things we changed
&
Things we would
do differently

04

LESSONS LEARNED

From
MongoDB
&
cloud data
management

05

APPENDIX

Dataset Size
&
Queries
&
GitHub Link



01

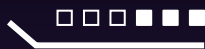
DATA



The Freeway Dataset, our MongoDB Data Model,
and our ETL Process



Freeway Dataset



- Relational data in .csv files
- Full of empty NULL values
- Querying in this format would require many joins
- Large freeway_loopdata file

detectorid	starttime	volume	speed	occupancy	status	dqflags
1345	9/15/2011 0:00:00		0		0	0
1345	9/15/2011 0:00:20				0	0
1345	9/15/2011 0:00:40	0			0	0
1345	9/15/2011 0:01:00	0			0	0
1345	9/15/2011 0:01:20	0			0	0
1345	9/15/2011 0:01:40	1	47		0	3
1345	9/15/2011 0:02:00	0			0	0
1345	9/15/2011 0:02:20	0			0	0

Table	A	Collection	{
Rows	of	Documents	_id: ObjectId(123),
Primary Key	identified by a	_id:ObjectId()	field: "value"
Columns	and contains		}
			{
			_id: ObjectId(124),
			field: "value"
			}

A NoSQL, distributed, document-oriented database that uses collections of JSON-type documents which support embedded fields for the storage of related data.



MongoDB

Detector Collection:

```
{
  _id:
    ObjectId("5fad9e5bfb3c57497c725fa5")
  detectorid:1345
  highway: {
    highwayid:3
    shortdirection:"N"
    direction:"NORTH"
    highwayname:"I-205"
  }
  milepost:14.32
  locationtext:"Sunnyside NB"
  detectorclass:1
  lanenumber:1
  station: {
    stationid:1045
    upstream:0
    downstream:1046
    stationclass:1
    numberlanes:4
    latlon:"45.43324,-122.565775"
    length:0.94
  }
}
```

OUR DATA MODEL

Freeway Loop Data Collection:

```
{
  _id: ObjectId("5fad9f01fb3c57497c725faf")
  detectorid:1345
  starttime:"2011-09-15 00:02:40-07"
  volume:1
  speed:66
}
```

EXTRACT TRANSFORM LOAD

```
1 #!/bin/bash
2 FNAME=freeway_loopdata.csv
3 HEADER=$(head -1 $FNAME)
4 split -b 300m $FNAME sections
5 n=1
6 for f in sections*
7 do
8     if [ $n -gt 1 ]; then
9         echo $HEADER > part${n}.csv
10    fi
11    cat $f >> part${n}.csv
12    rm $f
13    ((n++))
14 done
15
```

- The detectors, stations, and highways were in separate .csv files
- We wrote a Python script to format the data into a detector collection fitting our data model.
- This was saved to a JSON file and loaded to our database

- ← The raw freeway_loopdata.csv was too large to load into our database
- ← We used this split.sh script to section it into workable files
- ← We cleaned each file of NULL/0 values.
- ← This allowed us to load all data

```
1 import csv
2 import json
3 import pandas as pd
4
5 # read in csv files into dataframes
6 df_de = pd.read_csv ("data/freeway_detectors.csv")
7 df_st = pd.read_csv ("data/freeway_stations.csv")
8 df_hw = pd.read_csv ("data/highways.csv")
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65 #Save detector list to json file
66 with open("detectors.json", "w") as detectors_file:
67     json.dump(detectors, detectors_file, indent=4)
68
69
```



02

Queries

How we answered the freeway dataset questions



Connecting with Python

```
1 from pymongo import MongoClient
2 from pprint import pprint
3 import getpass as gp
4 pw = gp.getpass()
5 username = "DJs"
6 password = pw
7 dbname = "djs-freeway"
8 uri = "mongodb+srv://" + username + ":" + password + \
9       "@ccdm-project.f4c6t.mongodb.net/" + dbname + "?retryWrites=true&w=majority"
10 #client = MongoClient()
11 #client = MongoClient(uri)
12 #db = client.test
13 try:
14     client = MongoClient(uri)
15     db = client.test
16     print("Connected Successfully!!!")
17 except:
18     print("Could not connect to db :( ")
19
20
21 mydb = client[dbname]
22
23 de_collection = mydb["freeway_detectors"]
24 lp_collection = mydb["freeway_loopdata"]
25
```



Freeway Dataset Questions



Count low speeds and high speeds: Find the number of speeds < 5 mph and > 80 mph in the data set.



Volume: Find the total volume for the station Foster NB for Sept 15, 2011.



Single-Day Station Travel Times: Find travel time for station Foster NB for 5-minute intervals for Sept 15, 2011. Report travel time in seconds.



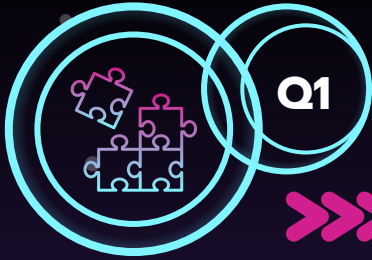
Peak Period Travel Times: Find the average travel time for 7-9AM and 4-6PM on September 22, 2011 for the I-205 NB freeway. Report travel time in minutes.



Route Finding: Find a route from Johnson Creek to Columbia Blvd on I-205 NB using the upstream and downstream fields.



Update: Change the milepost of the Foster NB station to 22.6.

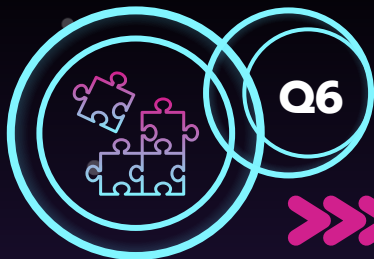


*Count low speeds and high speeds:
Find the number of speeds < 5 mph
and > 80 mph in the data set.*

```
# Query 1: count the number of speeds < 5 mph and > 80 mph
result1 = lp_collection.count_documents(
    {"speed": {"$lt": 5}})
print("Number of speeds < 5:", result1)
result2 = lp_collection.count_documents(
    {"speed": {"$gt": 80}})
print("Number of speeds > 80:", result2)
```

```
Number of speeds < 5: 1269204
Number of speeds > 80: 62203
```





Update: Change the milepost of the Foster NB station to 22.6.

```
32 #Query 6 update milepost at "Foster NB" from 18.1 -> 22.6
33 gry6 = de_collection.update_many({"locationtext": {"$eq": 'Foster NB'}}), {"$set": {"milepost": 22.6}})
34
35 cursor = de_collection.find({"locationtext": {"$eq": 'Foster NB'}})
36 for record in cursor:
37     print(record)
```

```
Matched:3
Updated:3
{'_id': ObjectId('5fc420ddce0db46a9c85b88d'), 'milepost': 18.1, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88e'), 'milepost': 18.1, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88f'), 'milepost': 18.1, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88d'), 'milepost': 22.6, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88e'), 'milepost': 22.6, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88f'), 'milepost': 22.6, 'locationtext': 'Foster NB'}
```



Clean up freeway loop data for the rest of the queries

```
84
85 # Remove unneeded loopdata after running Q1
86 # Q2 & 3 only need loopdata from Sept 15
87 # Q4 only needs loopdata from Sept 22
88 # Q5 onyl needs detector collection
89 filter={
90     '$nor': [
91         {
92             'starttime': {
93                 '$regex': '2011-09-22',
94                 '$options': 'i'
95             }
96         }, {
97             'starttime': {
98                 '$regex': '2011-09-15',
99                 '$options': 'i'
100             }
101         }
102     ]
103 }
104
105 result = lp_collection.delete_many(
106     filter=filter
107 )
108 print(f"Deleted:{result.deleted_count}")
109
```




03



SELF-CRITIQUE



Challenges we faced and how we would change our approach in the future.





QUERY PLANS

- First, used MongoDB Compass UI aggregation tool
- Switched to using a Python connection script
- After running Q1, we truncated freeway_loopdata due to "max space used (512mb) error stopping us from performing operations.

- Start figuring out queries earlier with smaller freeway_loopdata subset
- Start with Python connection script, then use Compass UI as needed for assistance



DATASET

[StackOverflow](#) (public dataset)

- Exported data from BigQuery to GCS bucket files
- Mounted GCS bucket to VM so we could access it like a local file system
- **AND THEN** 1 out of 1000s of files took up all 512MB of space in MongoDB

Switched to [Freeway Dataset](#)

- Do more research on dataset's size before trying to use it
- Find a subset that would work for StackOverflow dataset



ETL METHODS

- At first, we loaded in all four datasets as is, but realized we had to combine some of the datasets.
- Edited the data in the freeway_loopdata to contain the only the fields we needed for the queries.

- Design data model FIRST
- Write script to clean the freeway_loopdata of NULL/0 values
- Add another index, other than the default one if possible

DID CHANGE

WOULD
CHANGE



TIME INVESTMENT

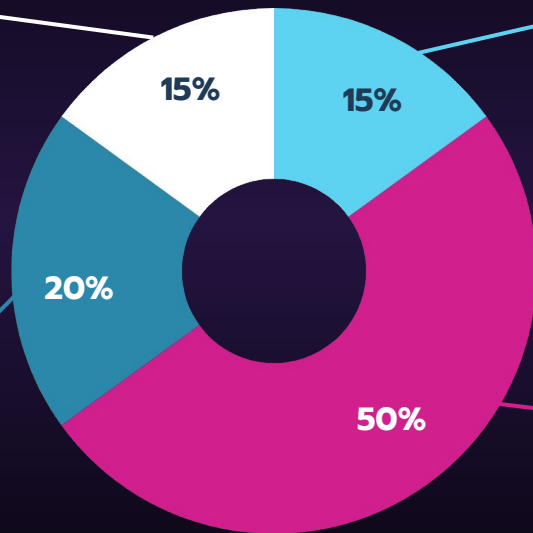


ADMIN

- Setting up MongoDB went pretty smoothly on our own boxes and we also set up a VM on GCP. We created a Github for scripts and a shared Google drive for cleaned and transformed data.

ARCHITECTURE

We spent a good amount of time thinking about the best data model with the goal of fitting it to the queries we needed to make.



ANALYZE

This has the potential to be a bit more than 15% by the time we complete them, but so far we have invested a disproportionate time doing ETL activities.

ETL

Choosing a dataset, cleaning and transforming the data, and loading it to our data model took a lot of time.



04



LESSONS LEARNED



MongoDB and Cloud Data Management



OUR TAKEAWAYS



- ETL takes time and planning
- MongoDB is a user-friendly database
- Dataset sizes can be very misleading
- Taking the time to write a script is better than trying to do the task manually
- MongoDB gives you a default index if you don't specify one
- Denormalized data > Normalized data

- Don't try to do ETL without deciding on a data model
- Some of the functions in the manual are deprecated, so check which versions supports the functions you're using.
- DB compass is great as a guide for writing some queries ...
- ... But just default to using a good connection script in a language you know



OUR ADVICE



QUESTIONS?

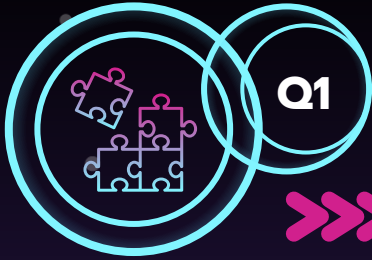
CREDITS: This presentation template was created by Slidesgo,
including icons by Flaticon, and infographics & images by Freepik

05

APPENDIX



Dataset Size: *Full Freeway Dataset (approximately 660MB)*
Number of Queries : 4



*Count low speeds and high speeds:
Find the number of speeds < 5 mph
and > 80 mph in the data set.*

```
# Query 1: count the number of speeds < 5 mph and > 80 mph
result1 = lp_collection.count_documents(
    {"speed": {"$lt": 5}})
print("Number of speeds < 5:", result1)
result2 = lp_collection.count_documents(
    {"speed": {"$gt": 80}})
print("Number of speeds > 80:", result2)
```

```
Number of speeds < 5: 1269204
Number of speeds > 80: 62203
```



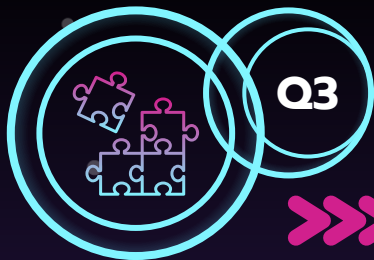


Volume: Find the total volume for the station Foster NB for Sept 15, 2011.

```
#Query 2 Find the total volume for the station Foster NB for Sept 15, 2011
qry2 = [{'$lookup': {'from': 'freeway_detectors',
                    'localField': 'detectorid',
                    'foreignField': 'detectorid',
                    'as': 'detectors'}},
        {'$match': {'detectors.locationtext': 'Foster NB'}},
        {'$match': {'starttime': {'$regex': '2011-09-15'}}},
        {'$group': {'_id': 'None', 'TotalVolume': {
                    '$sum': '$volume'}}}]
cursor = lp_collection.aggregate(qry2)
result = list(cursor)
print("Query 2 results: ",result)
```

```
Query 2 results:  [{'_id': 'None', 'TotalVolume': 49891}]
```





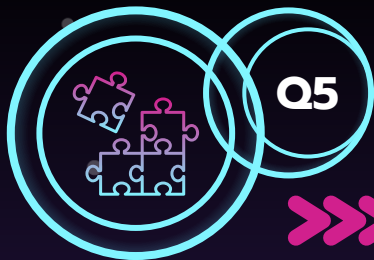
Single-Day Station Travel Times: Find travel time for station Foster NB for 5-minute intervals for Sept 15, 2011. Report travel time in seconds.





Peak Period Travel Times: Find the average travel time for 7-9AM and 4-6PM on September 22, 2011 for the I-205 NB freeway. Report travel time in minutes.





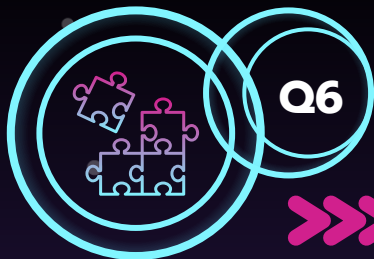
Route Finding: Find a route from Johnson Creek to Columbia Blvd on I-205 NB using the upstream and downstream fields.

```
#Query 5 Find the path from Johnson Creek to I-205 NB at Columbia
```

```
i = 0
text = 'Johnson Cr NB'
print(i,":",text)
while text != 'I-205 NB at Columbia':
    qry5 = [ {'$match': {'locationtext': text}},
             {'$lookup': {
                 'from': 'freeway_detectors',
                 'let': {'down': '$station.downstream',
                        'lanenum': '$lanenum'},
                 'pipeline': [ {'$match': {'$expr': {
                                '$eq': ['$station.stationid', '$$down']}},
                              {'$match': {'$expr': {'$eq': ['$lanenum', '$$lanenum']}}}
                             ],
                 'as': 'downstation'},
             {'$limit': 1},
             {'$unwind': {'path': '$downstation'}},
             {'$project': {'downstation.locationtext': 1}}]
    cursor = de_collection.aggregate(qry5)
    result = list(cursor)
    i += 1
    for doc in result:
        text2 = doc["downstation"]
        text = text2["locationtext"]
        print(i,":",text)
```

```
0 : Johnson Cr NB
1 : Foster NB
2 : Powell to I-205 NB
3 : Division NB
4 : I-205 NB at Glisan
5 : I-205 NB at Columbia
```





Update: Change the milepost of the Foster NB station to 22.6.

```
32 #Query 6 update milepost at "Foster NB" from 18.1 -> 22.6
33 gry6 = de_collection.update_many({"locationtext": {"$eq": 'Foster NB'}}), {"$set": {"milepost": 22.6}})
34
35 cursor = de_collection.find({"locationtext": {"$eq": 'Foster NB'}})
36 for record in cursor:
37     print(record)
```

```
Matched:3
Updated:3
{'_id': ObjectId('5fc420ddce0db46a9c85b88d'), 'milepost': 18.1, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88e'), 'milepost': 18.1, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88f'), 'milepost': 18.1, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88d'), 'milepost': 22.6, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88e'), 'milepost': 22.6, 'locationtext': 'Foster NB'}
{'_id': ObjectId('5fc420ddce0db46a9c85b88f'), 'milepost': 22.6, 'locationtext': 'Foster NB'}
```





Source code available on Github

<https://github.com/santitobon9/Cloud-Cluster-Freeway-Project.git>

