STAT 405 Final Presentation

By Ella Gruen, Anthony Pagas, John Chumlea, Shien Zhu and Samuel Negus

Introduction

Dataset: NYC Taxis

- Taxi trip records exclusively in NYC with variables such as pickup + dropoff time, trip details, passenger count, and vendor ID

Research Question: Which taxi vendor is the most time efficient to take during any day of the week?

Statistical Computation Method: **Group-wise aggregation**, we grouped trip data by vendor ID and day of the week and calculated the mean of minutes per mile within each group. Ex: "On Tuesday around 12 AM, CMT taxis traveling between 0-1 mile averaged about 5.39 minutes per mile"

Data Example

	vendor_id	distance_bin	hour	dayofweek	minutes_per_mile
1	СМТ	0–1mi	0	0	5.61716386305714 7
2	СМТ	0–1mi	0	1	5.3848840462477 32
3	СМТ	0–1mi	0	2	5.56614460175915 6
4	СМТ	0–1mi	0	3	5.5839515297479 51
5	СМТ	0–1mi	0	4	6.49177598027821 2
6	СМТ	0–1mi	0	5	7.03520771089674 3
7	СМТ	0–1mi	0	6	7.65842897101342 1
8	СМТ	0–1mi	1	0	4.9814533139552 57
9	СМТ	0–1mi	1	1	5.2683728820742 94
10	СМТ	0–1mi	1	2	5.57912941783701 9

Data

Source/Size: Kaggle (<u>https://www.kaggle.com/datasets/chilam/nyctaxis</u>), 28.85 GB

Cleaning:

- Dropped all trip distances that were not a valid number
- Converted "pickup_datetime" variable into a datetime object and dropped any corrupt time formats
- Removed trips with zero or negative distance or time

Computation:

- We ran 12 parallel jobs using 2 GB of request disk/memory and 1 CPU. All jobs lasted about 10-20 minutes

Key Points of Coding

- efficiency.py
 - 2GB/ csv, so doing computation for each creates an issue in the loop.
 - Idea: Using chunk library, we can compute 100,000 rows each iteration.
- Parallel job
 - Created a python environment, similar to the way we were able to run Python scripts in HW4.
- Issues:
 - The .csv file is so large that the compute node has a difficult time doing computation -> change python library (panda -> chunk)
 - HTC doesn't have Python environment to run our code -> Create Python environment, run code in the environment



Overall Vendor Efficiency (Lower = Faster)

		vendor_id	minutes_per_mile
1	1	VTS	5.1159957909080 6
2	0	СМТ	5.6303473699478 03

Vendor Efficiency By Day Of Week

	vendor_id	dayofweek	minutes_per_mile
1	СМТ	0	4.3992049868345 61
2	СМТ	1	4.6536560653587 34
3	СМТ	2	4.7410015250485 89
4	СМТ	3	4.7720092672802 66
5	СМТ	4	4.8389871503246 9
6	СМТ	5	5.0623425967464 03
7	СМТ	6	10.945229998041 37
8	VTS	0	4.8898203067201 695
9	VTS	1	5.2799605433133 34
10	VTS	2	5.3190167560585 05
11	VTS	3	5.2630047292688 68
12	VTS	4	5.7956703578513 97
13	VTS	5	4.8402274974283 52
14	VTS	6	4.4242703457157 9

Vendor Efficiency By Distance Bin

	vendor_id	distance_bin	minutes_per_mile
1	СМТ	0–1mi	10.1977337850414 86
2	СМТ	10mi+	2.7193125560845 606
3	СМТ	1–3mi	6.4943755709677 18
4	СМТ	3–6mi	5.1584939856721 97
5	СМТ	6–10mi	3.5818209519730 47
6	VTS	0–1mi	10.206902408405 36
7	VTS	10mi+	2.3339834471347 545
8	VTS	1–3mi	5.6702720546749 6
9	VTS	3–6mi	4.35381253271142 2
10	VTS	6–10mi	3.01500851161379 9

Graphical Representations









Conclusion

- VTS is the most efficient vendor in our NYC Taxi Data, averaging about 5 minutes and 6 seconds per mile
- Even though CMT vendor beats out VTS in minutes per mile 5 out of the 7 days of the week, VTS beats CMT in every distance bin category except for 0-1 mile
- We struggled a lot with disk space and python processing
- Prevent data inconsistencies by implementing distance bin (ex: if CMT had the second fastest time of day but only had one trip but VTS had the fastest time of the day but had multiple trips, CMT has a bias)