

Analysis of NYC Taxi & Limousine Commission Dataset (2012 - 2021)

Alec Peterson
DSCI 521-900
Summer 2022

A photograph of a dense traffic jam of yellow taxis in a city street. The taxis are packed closely together, filling the frame. The text "Project Overview" is overlaid in the center in a large, white, sans-serif font. The background is slightly blurred, emphasizing the taxis in the foreground.

Project Overview

polars for “big” data Processing

- To experiment and process even a year’s worth of data (~9 GB), used the `polars` module
- `polars` is a columnar query engine written in the low-level and increasingly popular **Rust** programming language, and is also made available as a Python package
- Has a similar syntax to `pandas`, but is significantly more performant for large datasets due to:
 - Low-level optimizations
 - Parallelism (using more than one CPU core)
 - Optimized in-memory columnar representation via Apache Arrow
- `polars` worked better out-of-the-box than other packages like `dask` (which parallelizes python code) or `pyarrow` (which deals more directly with Arrow structures)



Memory limits still a challenge

- Combined Taxi and For-Hire Vehicle data from 2012 through 2021 would have exceeded computer's RAM limits (32 GB)
 - ➔ Had to experiment on smaller datasets (1 mo – 1 yr)
- **Process:**
 1. Experiment with `polars` methods on small dataset
 2. Convert groupby-aggregated dataset (in MB) to `pandas` or write to more manageable `.parquet` file for subsequent `pandas` analysis
 3. Graph with `pandas`-compatible libraries (especially `plotly`)

plotly.express

- `plotly` is an open-source JavaScript library with many interactive chart types, with many made easily available in Python with `plotly.express`
- Python API straightforward
 - Less clunky and more visually appealing than `matplotlib`
 - Like `seaborn` but with more customization and interactivity
- Professional-looking plots, and interactivity valuable for dashboarding applications or other reporting
- Includes charts for Maps, like choropleths which were best for visualizing NYC Taxi Zones (simpler than using `matplotlib` with `geopandas`)

Data Analysis

A close-up photograph of a yellow taxi sign with the word "TAXI" in black capital letters. The sign is mounted on a chrome pole. The background is a blurred city street at night, with various colorful lights (red, yellow, blue, green) creating a bokeh effect. The text "Data Analysis" is overlaid in white, sans-serif font across the center of the image.

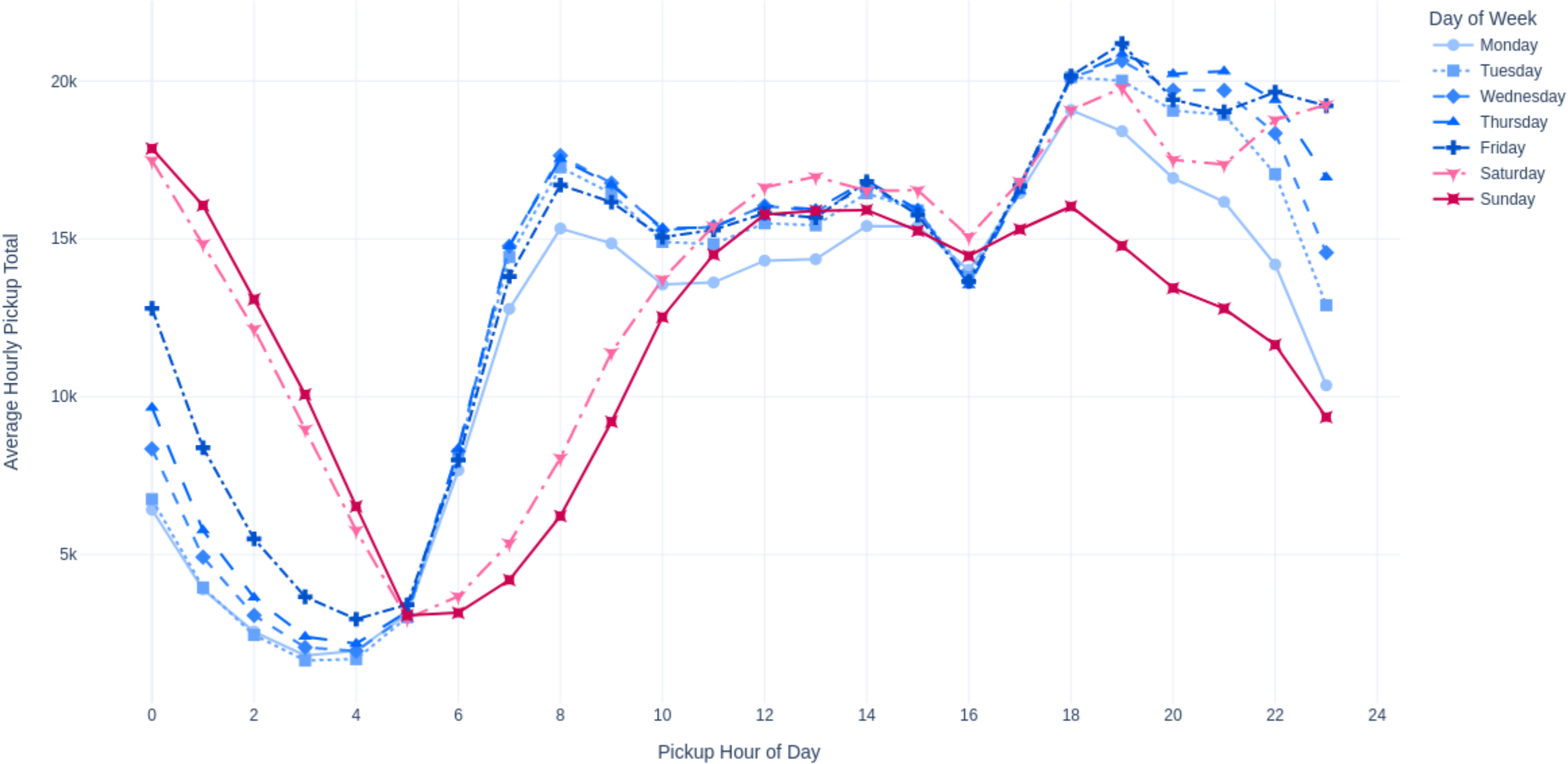
Analysis Goals

1. Temporal Analyses

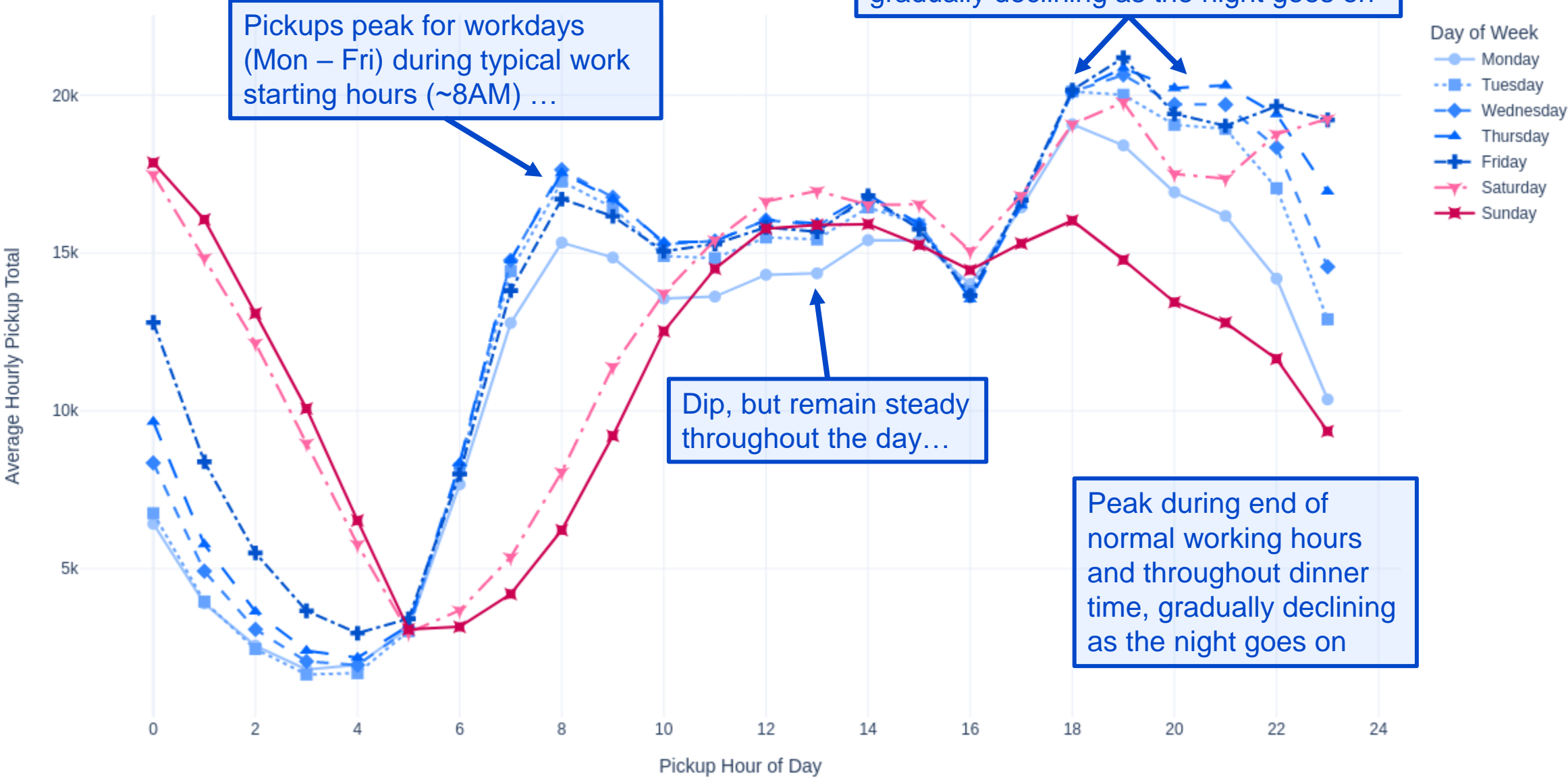
- A. “Busy” periods over multiple time scales (week, year, multiple years)
- B. Impact of for-hire vehicles (e.g. Uber and Lyft)
- C. Impact of the COVID-19 pandemic

2. Geospatial Analysis - busiest areas/boroughs (also over time)

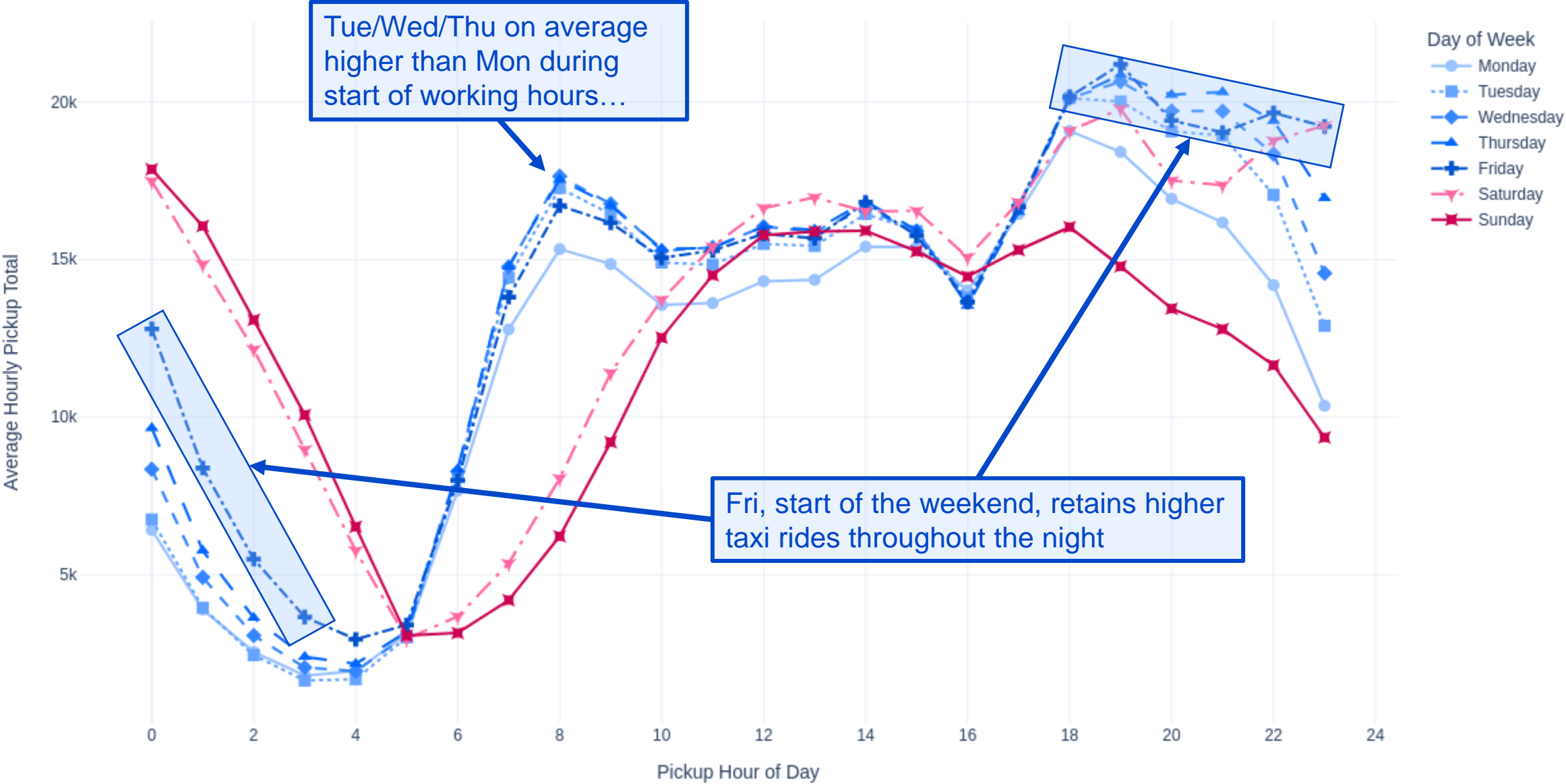
Average Hourly Pickup Total vs. Hour of Day, by Day of Week (Yellow Taxi)



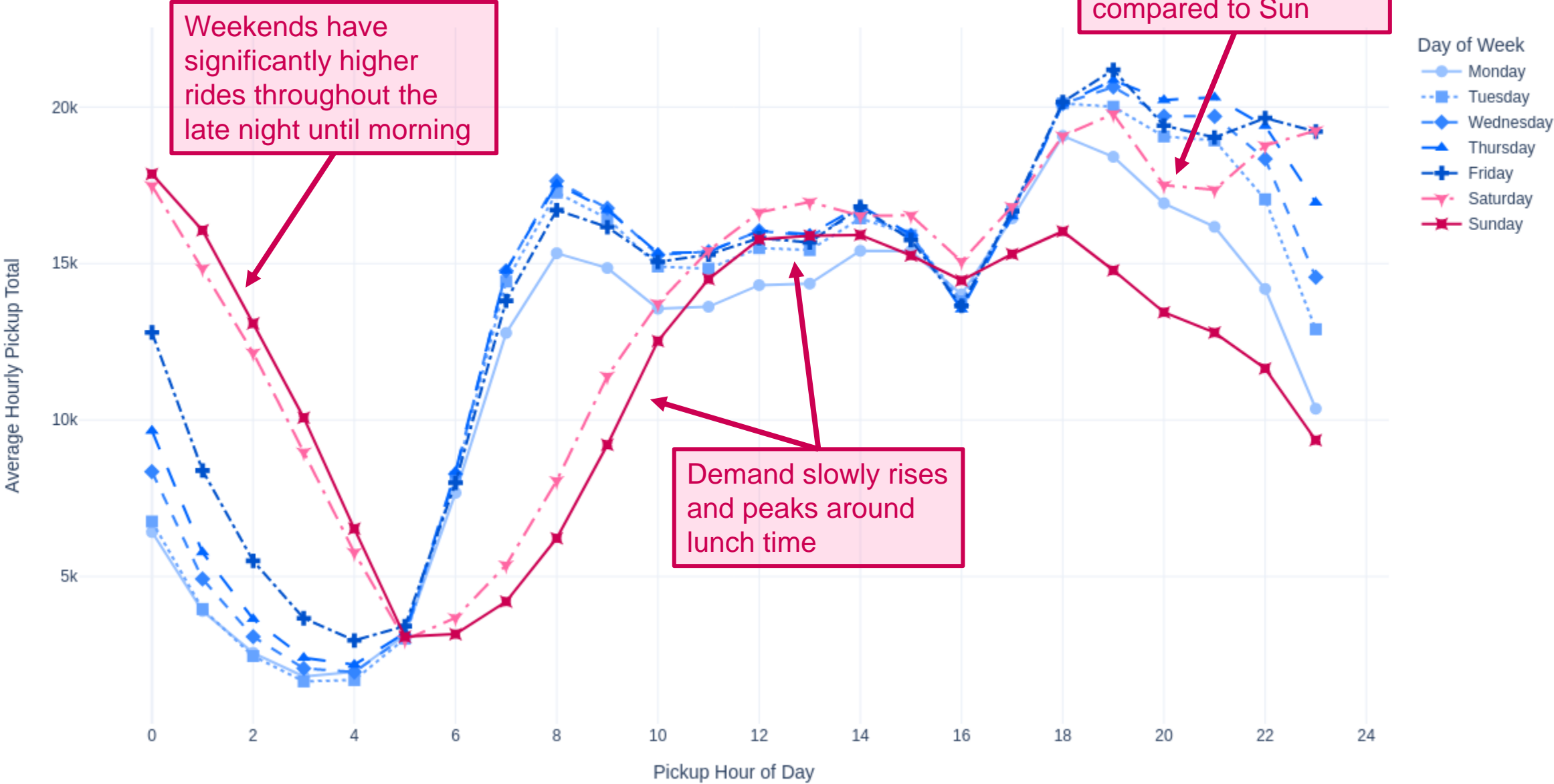
Average Hourly Pickup Total vs. Hour of Day, by Day of Week (Yellow Taxi)



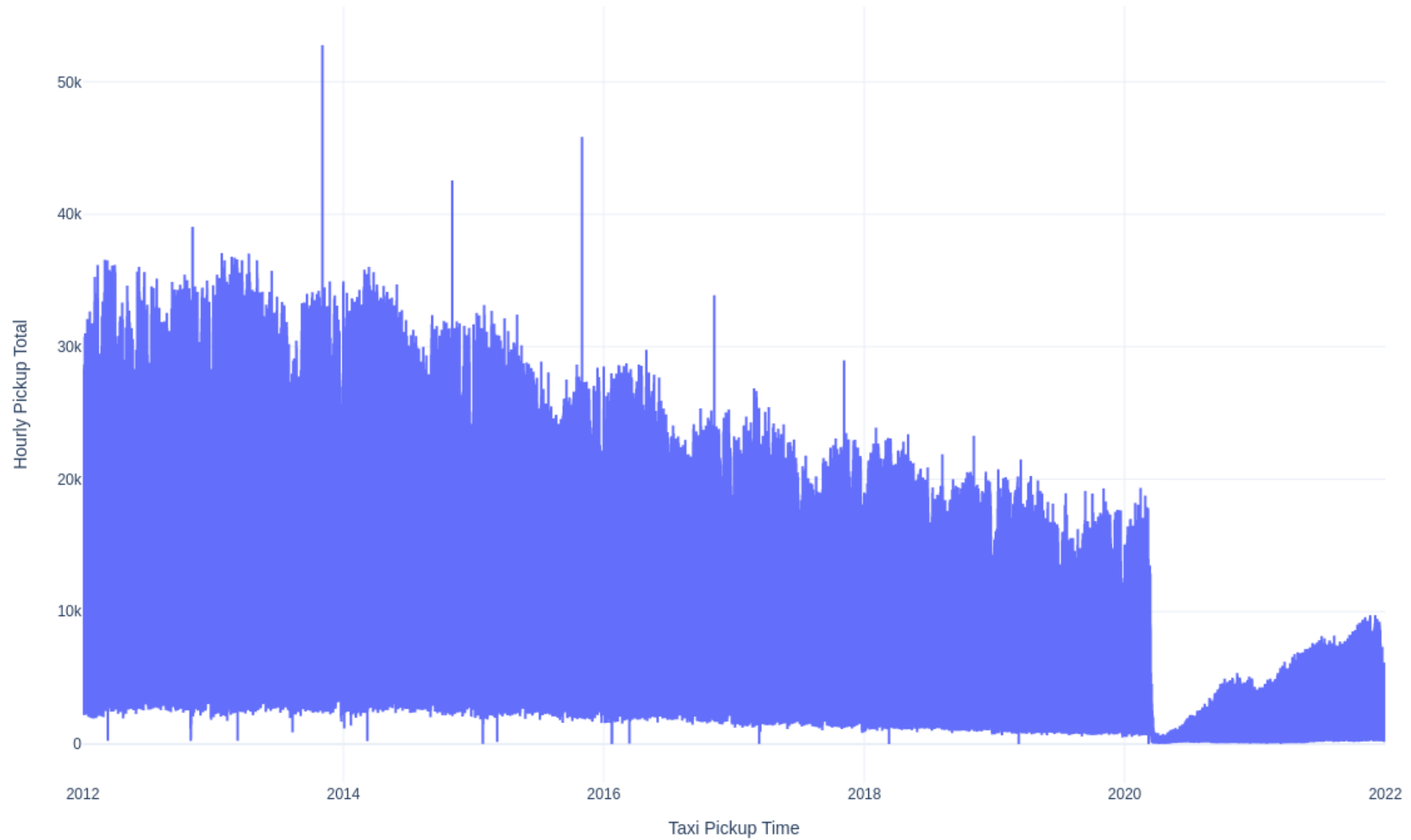
Average Hourly Pickup Total vs. Hour of Day, by Day of Week (Yellow Taxi)



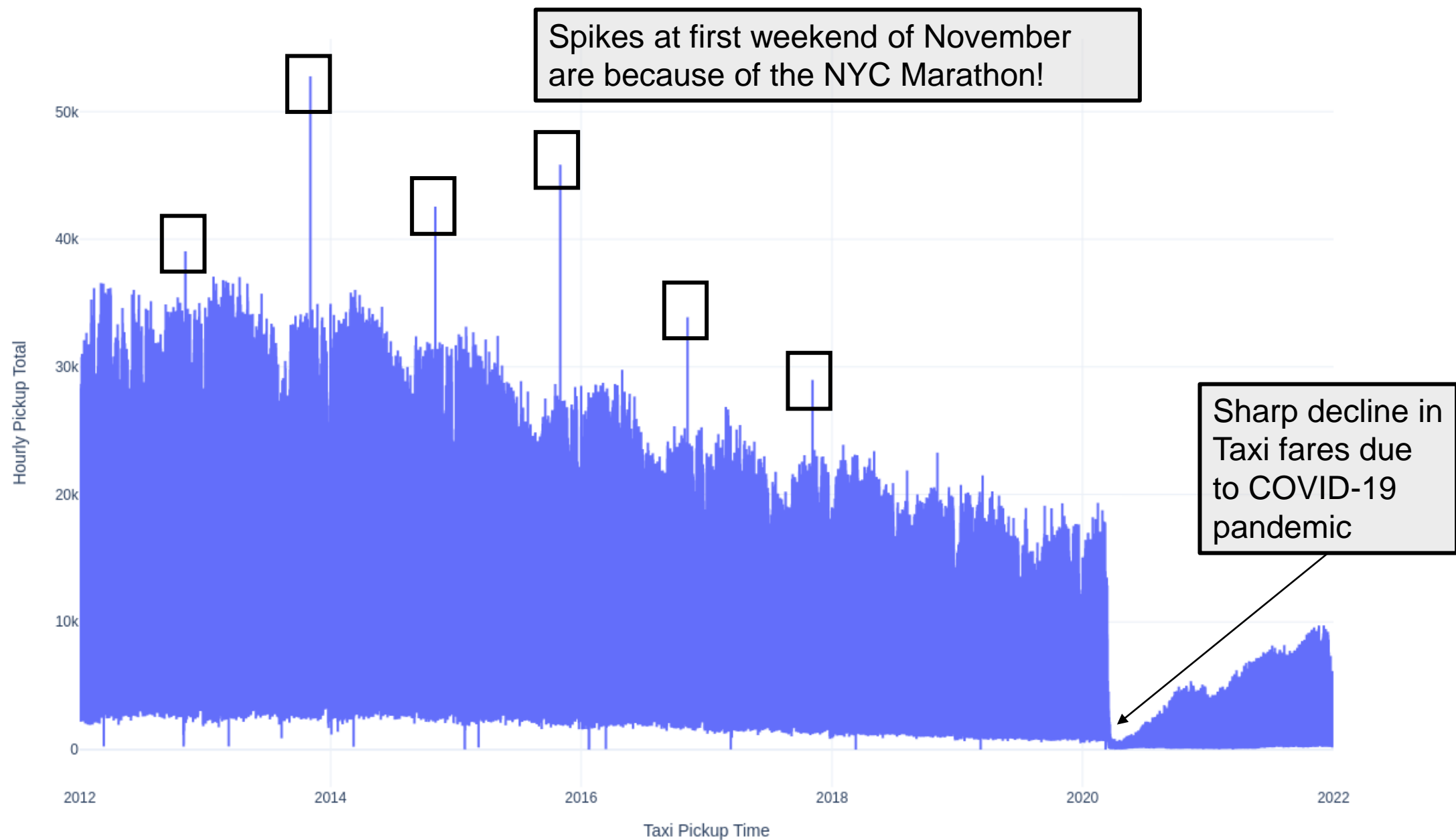
Average Hourly Pickup Total vs. Hour of Day, by Day of Week (Yellow Taxi)



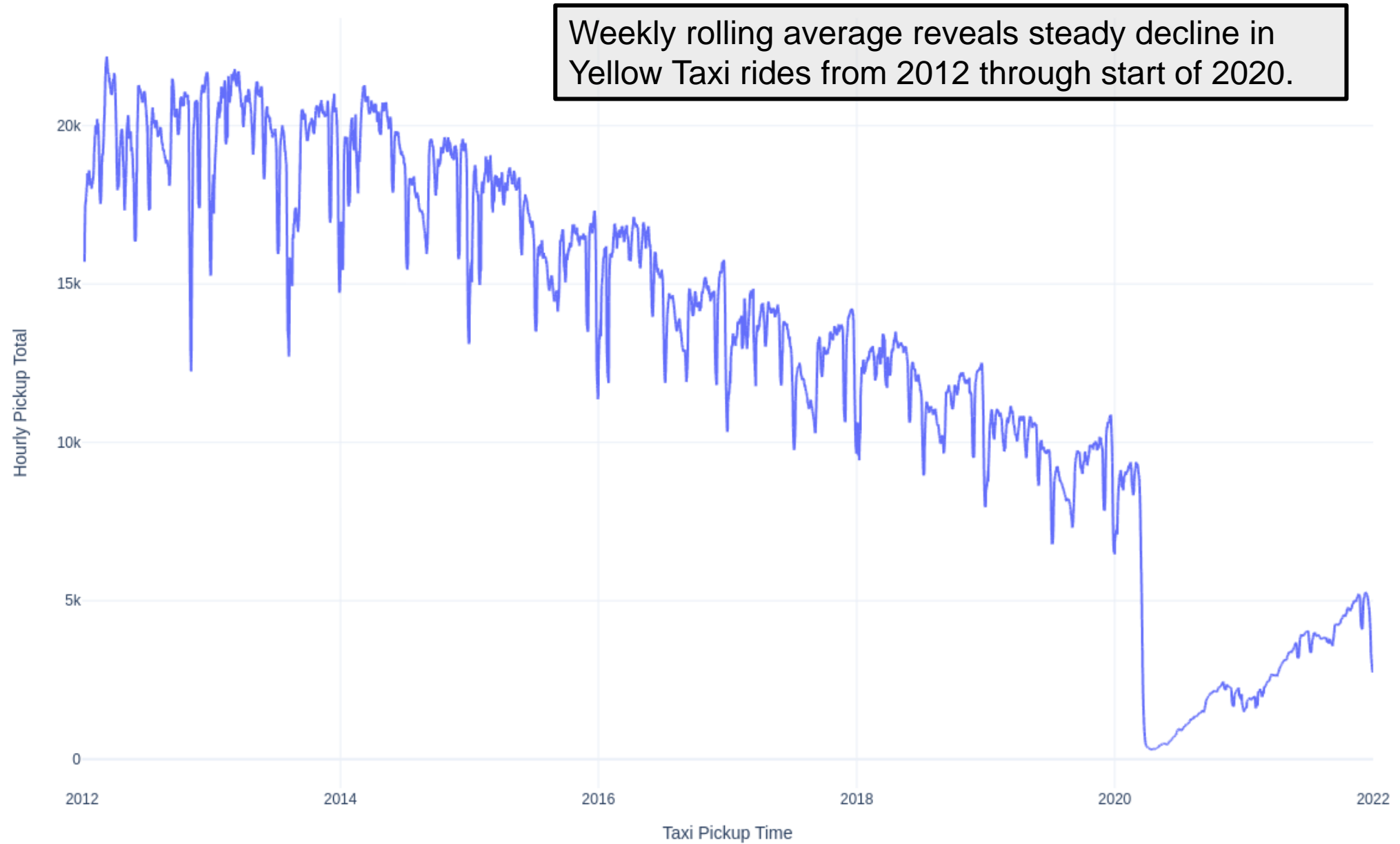
Hourly Pickup Total vs. Time



Hourly Pickup Total vs. Time



Hourly Pickup Total vs. Time



Analysis Goals

1. Temporal Analyses

A. “Busy” periods over multiple time scales (week, year, multiple years)

B. Impact of for-hire vehicles (e.g. Uber and Lyft)

C. Impact of the COVID-19 pandemic

2. Geospatial Analysis - busiest areas/boroughs (also over time)

Green Cabs and For-Hire Vehicles

- In August 2013, apple-green painted taxis (“**Green Cabs**”) were allowed to pick up passengers in Upper Manhattan, the Bronx, Brooklyn, Queens (excluding LaGuardia and JFK Airports) and Staten Island.
- In ~2015, **For-Hire Vehicles (FHVs)** were allowed to provide pre-arranged rides (e.g. via ride-hailing apps).
- In 2019, a “**High Volume**” category of FHVs was created for dispatch bases with >10,000 trips per day. This largely represents cars from the most popular ride-hailing services, Uber and Lyft.



Uber

lyft

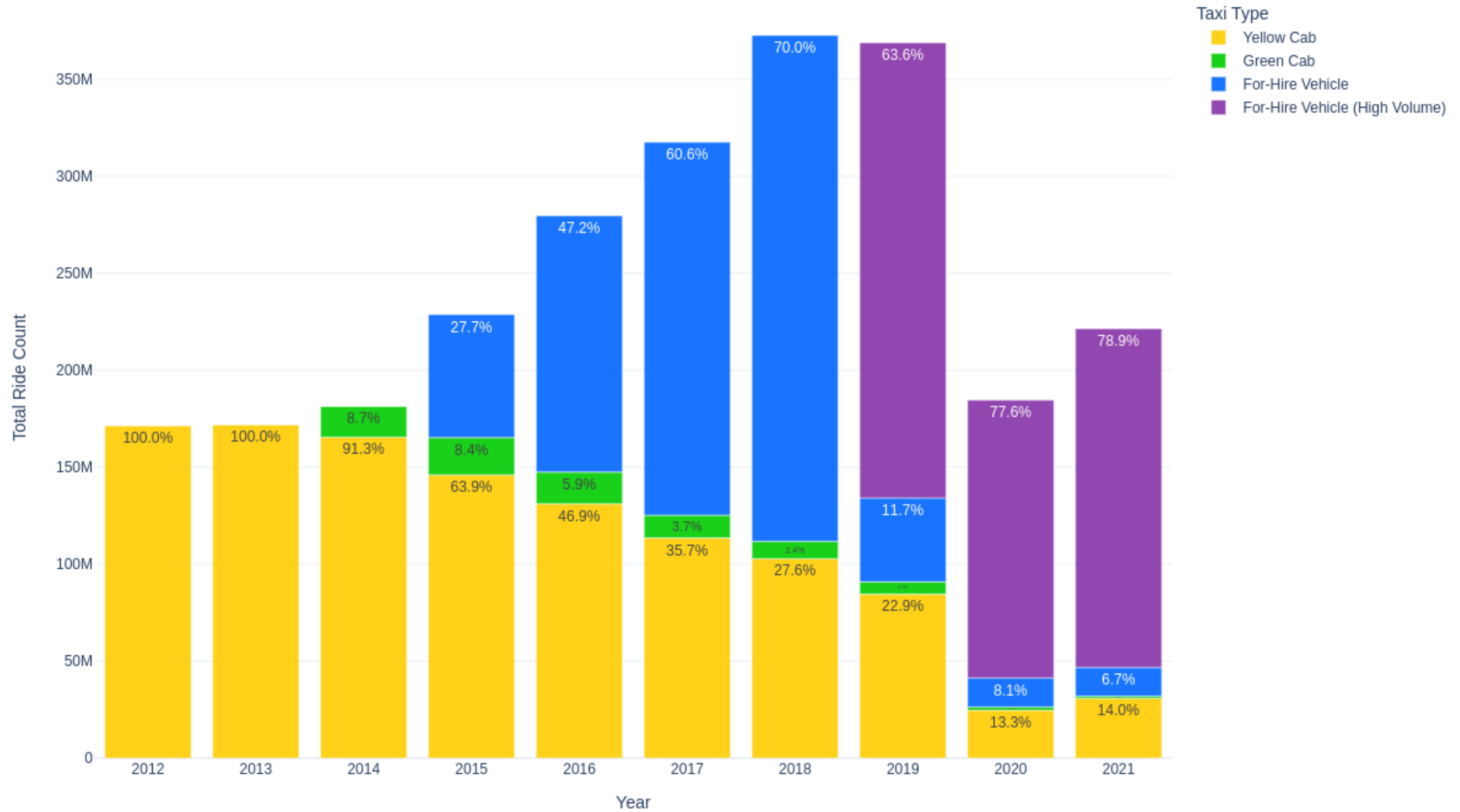
Hourly Pickup Total vs. Time



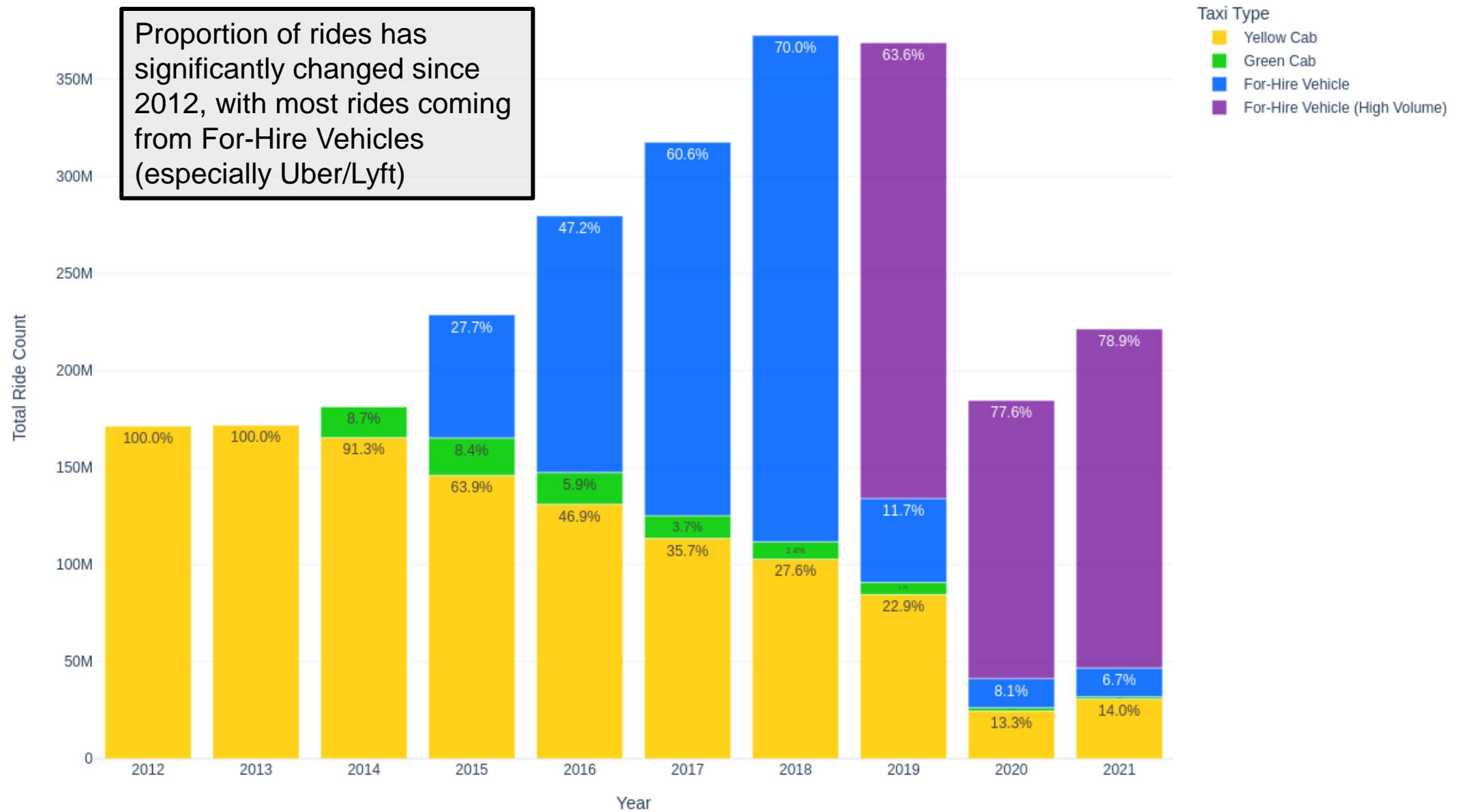
Hourly Pickup Total vs. Time



Yearly Total Ride Count vs. Time (years), by Taxi Type



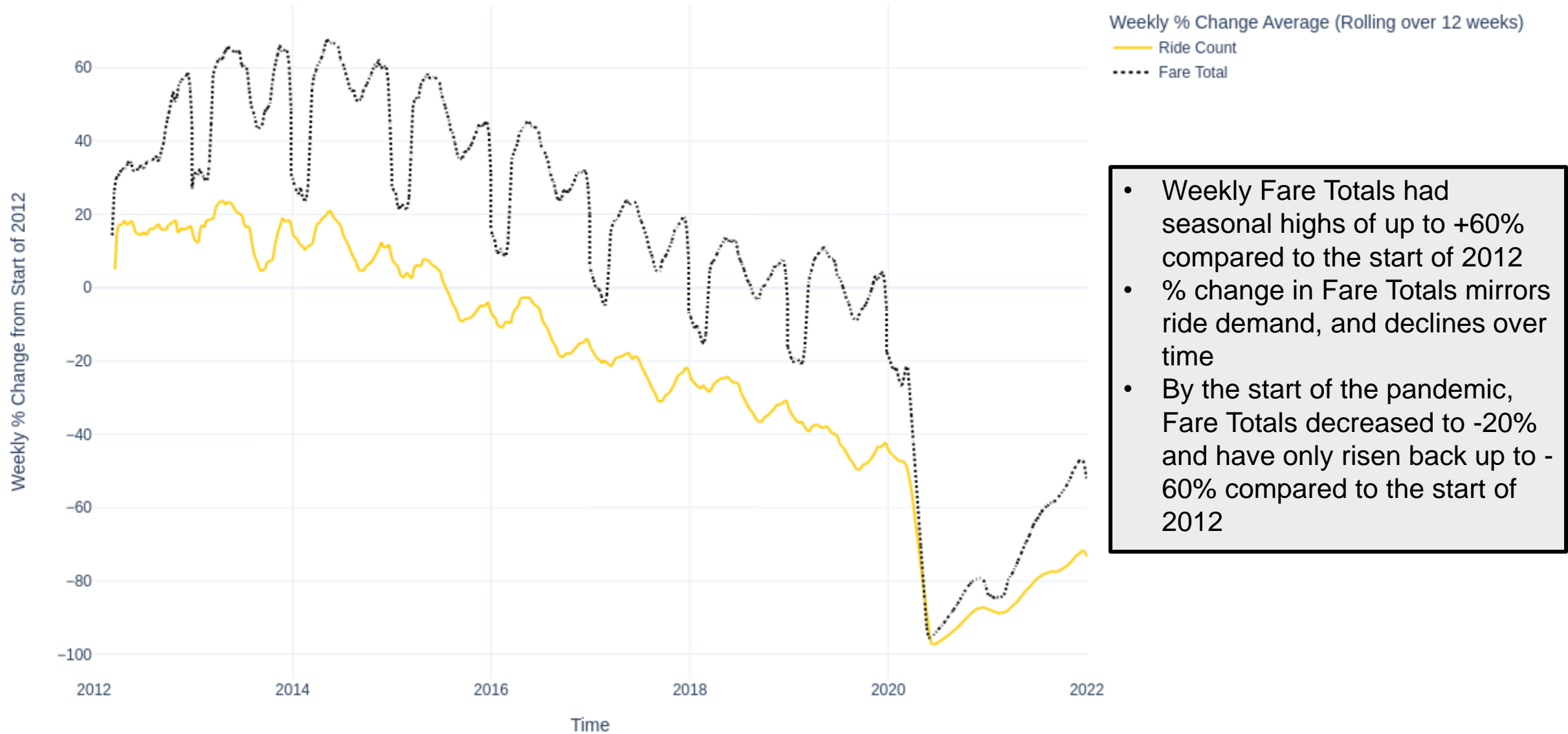
Yearly Total Ride Count vs. Time (years), by Taxi Type



Weekly Total Fare vs. Time for Yellow Taxis



Weekly Ride Count and Fare Total % Change (compared to 2012 Start) vs. Time



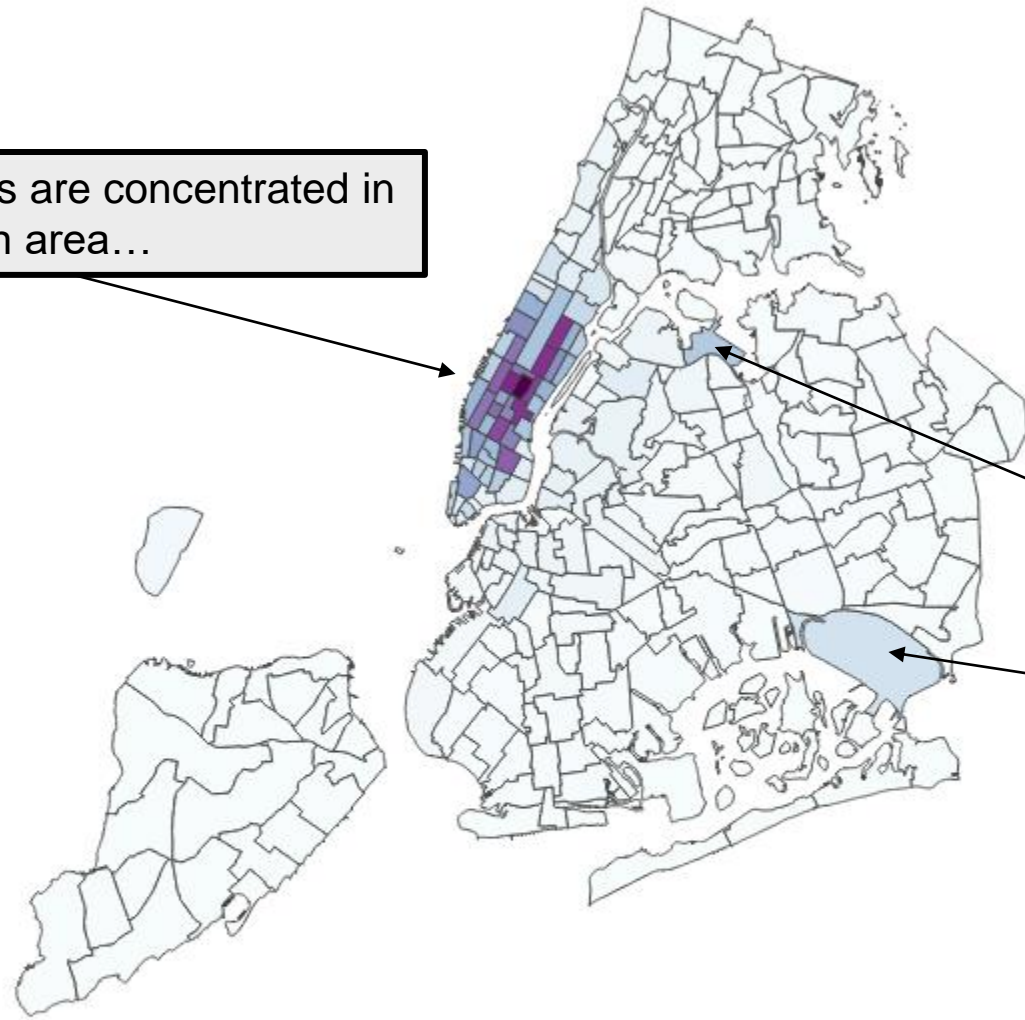
Analysis Goals

1. Temporal Analyses

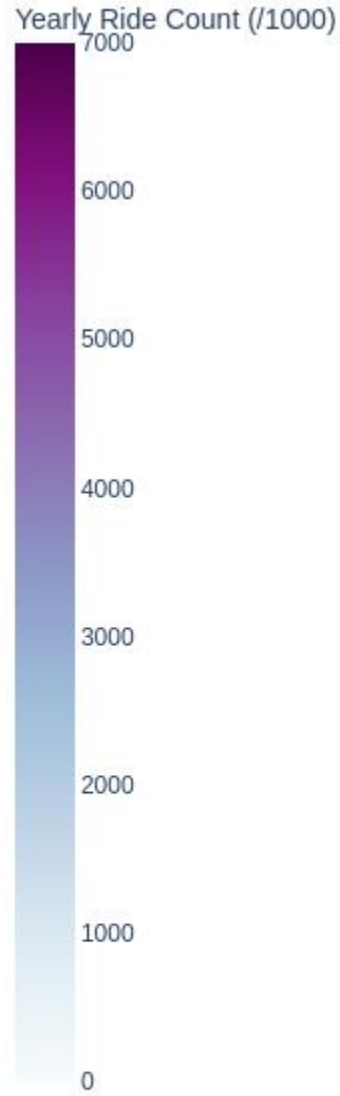
- A. “Busy” periods over multiple time scales (week, year, multiple years)
- B. Impact of for-hire vehicles (e.g. Uber and Lyft)
- C. Impact of the COVID-19 pandemic

2. Geospatial Analysis - busiest areas/boroughs (also over time)

Most taxi fares are concentrated in the Manhattan area...



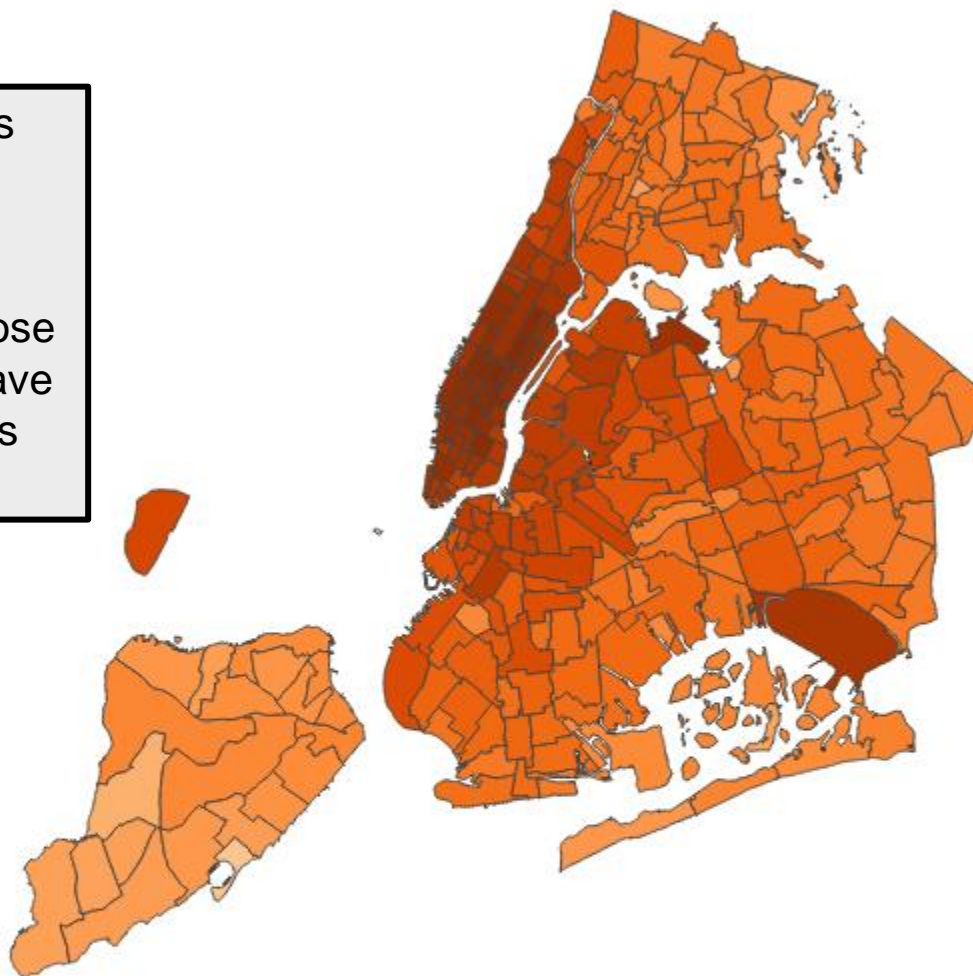
...as well as LaGuardia and JFK airports.



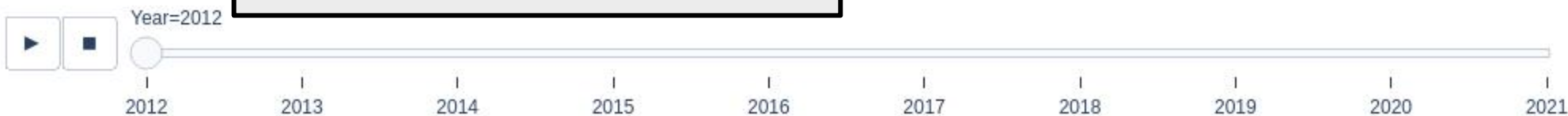
Jupyter Notebook interactive `plotly` version can show evolution over time!



- A log analysis of taxi zones gives a better sense of magnitude of rides throughout NYC
- Other areas, especially those close to Manhattan, still have significant Yellow Cab rides per year



Jupyter Notebook interactive `plotly` version can show evolution over time!



Conclusions

1. Takeaways

- Yellow taxi industry significantly impacted by ride-hailing apps, representing only 14% of total rides in 2021 compared to 100% in 2012
- COVID-19 pandemic significantly decreased
- Temporal and geospatial analysis gives insight to periods and areas of peak demand

2. Future Analysis Directions

- Analysis of airport taxi pickups/drop-offs throughout the day and proportion of total rides in NYC
- Analysis of busiest taxi zones within each borough
- Parallel computing to work with dataset more directly (multiple computers), GPU acceleration
- Slides with revealjs to embed visualizations within slides

Analysis of NYC Taxi & Limousine Commission Dataset (2012 - 2021)

Alec Peterson
DSCI 521-900
Summer 2022