

# Stream Processing

Project presentation

2020/2021

# Table of Contents

- Data - Taxi trips from New York City
- Questions
- Recommendations

# Data - Taxi trips from New York City

# Data

- Provided data consists of reports of **taxis trips** including **starting point**, **drop-off point**, corresponding **timestamps**, and information related to the **payment**.
- Data are reported at the end of the trip, i.e., upon arrive in the order of the drop-off timestamps.
- Events with the same `dropoff_datetime` are in random order.
- Quality of the data is not perfect.
  - Some events might miss information such as drop off and pickup coordinates or fare information.
  - Moreover, some information, such as, e.g., the fare price might have been entered incorrectly by the taxi drivers thus introducing additional skew.

# Taxi trips from New York City

Attributes	Description
<b>medallion</b>	an md5sum of the identifier of the taxi - vehicle bound
<b>hack_license</b>	an md5sum of the identifier for the taxi license
<b>pickup_datetime</b>	time when the passenger(s) were picked up
<b>dropoff_datetime</b>	time when the passenger(s) were dropped off
<b>trip_time_in_secs</b>	duration of the trip
<b>trip_distance</b>	trip distance in miles

# Taxi trips from New York City

Attributes	Description
<b>pickup_longitude</b>	longitude coordinate of the pickup location
<b>pickup_latitude</b>	latitude coordinate of the pickup location
<b>dropoff_longitude</b>	longitude coordinate of the drop-off location
<b>dropoff_latitude</b>	latitude coordinate of the drop-off location

# Taxi trips from New York City

Attributes	Description
<b>payment_type</b>	the payment method - credit card or cash
<b>fare_amount</b>	fare amount in dollars
<b>surcharge</b>	surcharge in dollars
<b>mta_tax</b>	tax in dollars
<b>tip_amount</b>	tip in dollars
<b>tolls_amount</b>	bridge and tunnel tolls in dollars
<b>total_amount</b>	total paid amount in dollars

# Taxi trips from New York City

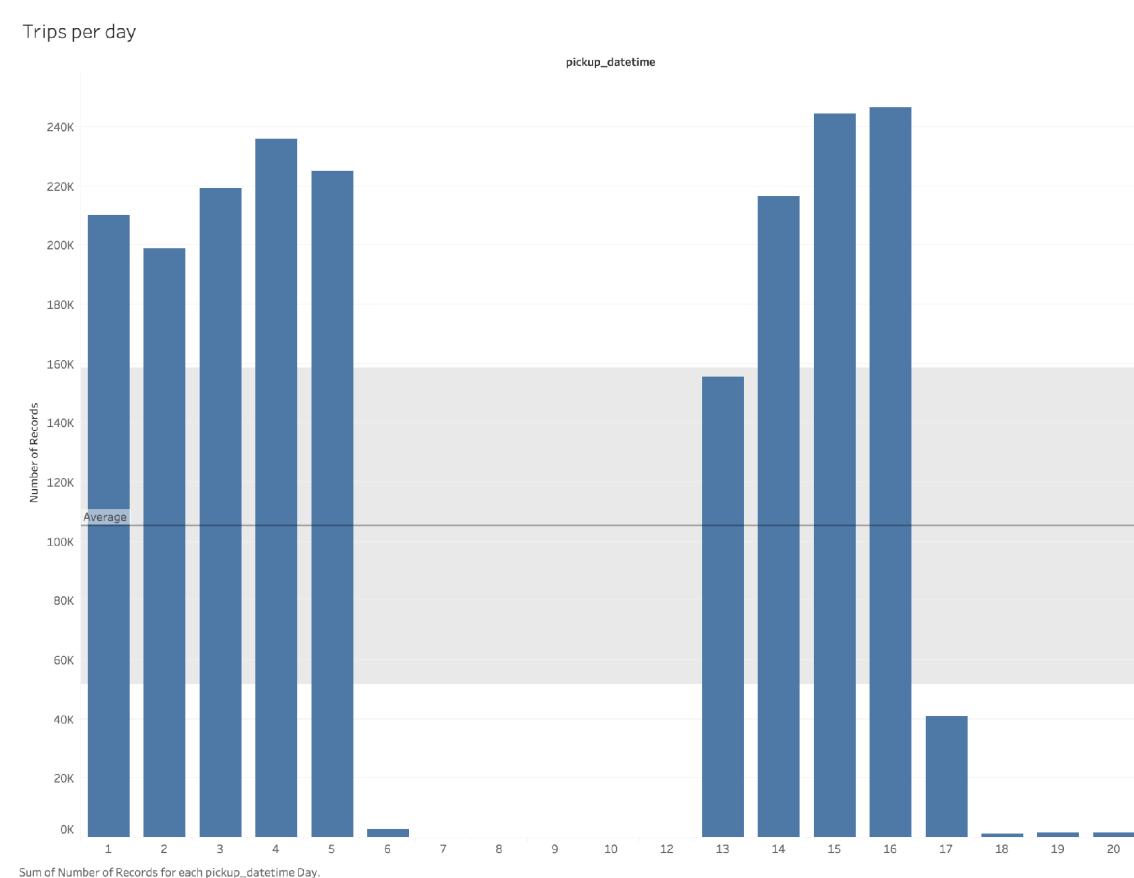
<b>medallion</b>	an md5sum of the identifier of the taxi - vehicle bound
<b>hack_license</b>	an md5sum of the identifier for the taxi license
<b>pickup_datetime</b>	time when the passenger(s) were picked up
<b>dropoff_datetime</b>	time when the passenger(s) were dropped off
<b>trip_time_in_secs</b>	duration of the trip
<b>trip_distance</b>	trip distance in miles
<b>pickup_longitude</b>	longitude coordinate of the pickup location
<b>pickup_latitude</b>	latitude coordinate of the pickup location
<b>dropoff_longitude</b>	longitude coordinate of the drop-off location
<b>dropoff_latitude</b>	latitude coordinate of the drop-off location
<b>payment_type</b>	the payment method - credit card or cash
<b>fare_amount</b>	fare amount in dollars
<b>surcharge</b>	surcharge in dollars
<b>mta_tax</b>	tax in dollars
<b>tip_amount</b>	tip in dollars
<b>tolls_amount</b>	bridge and tunnel tolls in dollars
<b>total_amount</b>	total paid amount in dollars

# Where to get the data?

- ACM DEBS 2015 Grand Challenge
  - <http://www.debs2015.org/call-grand-challenge.html>
- 20 days (roughly 2 million events) of data (~130 MB)
- Data for the whole year 2013 (~173 million events) (~12 G) (~33 G expanded)

# Preliminary information based on sample data

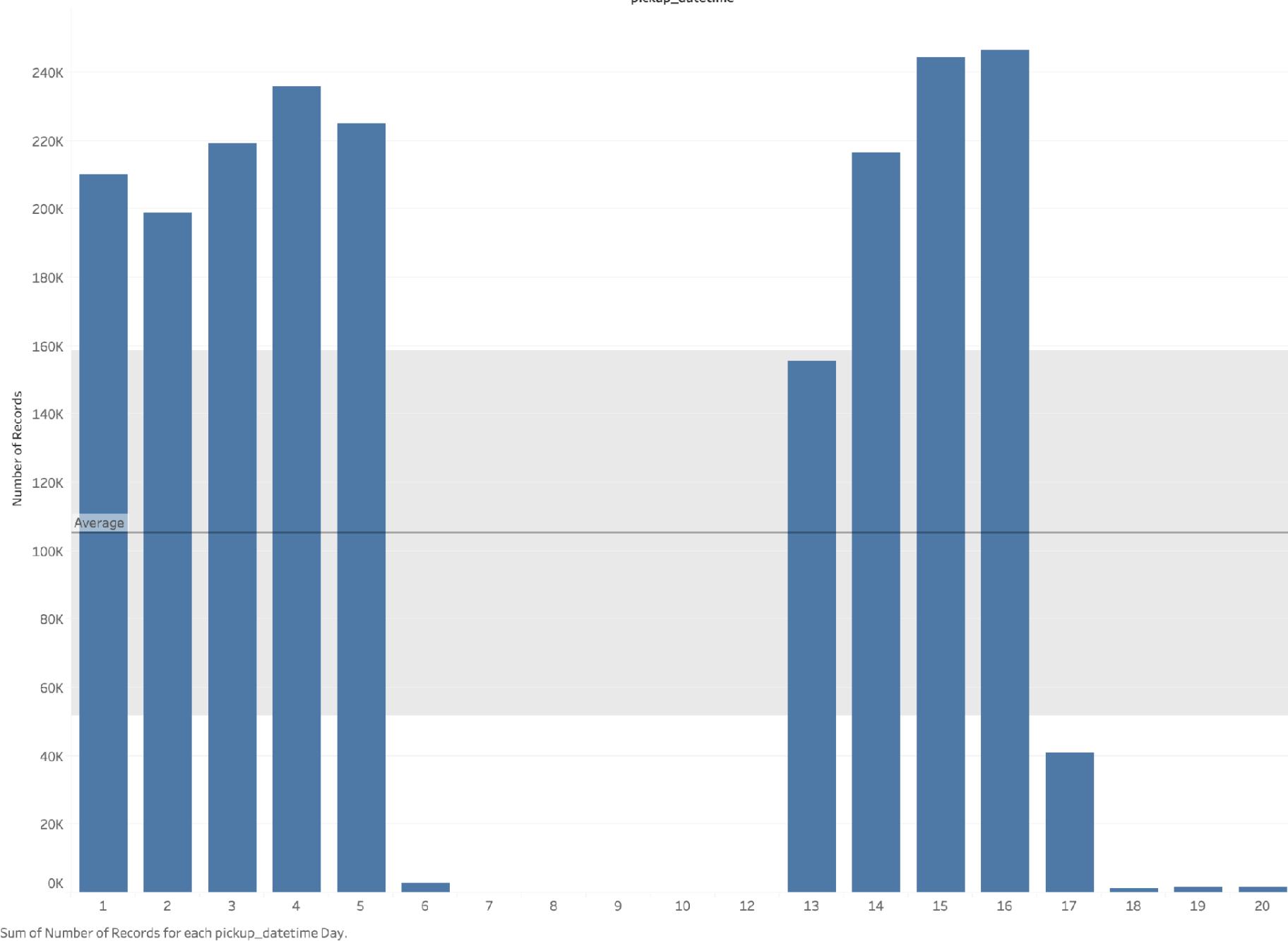
- ~10 800 Taxis
- ~20 300 Drivers
- 20 days
- ~2 million records (trips)



Trips per day

## Distribution of number of trips per day

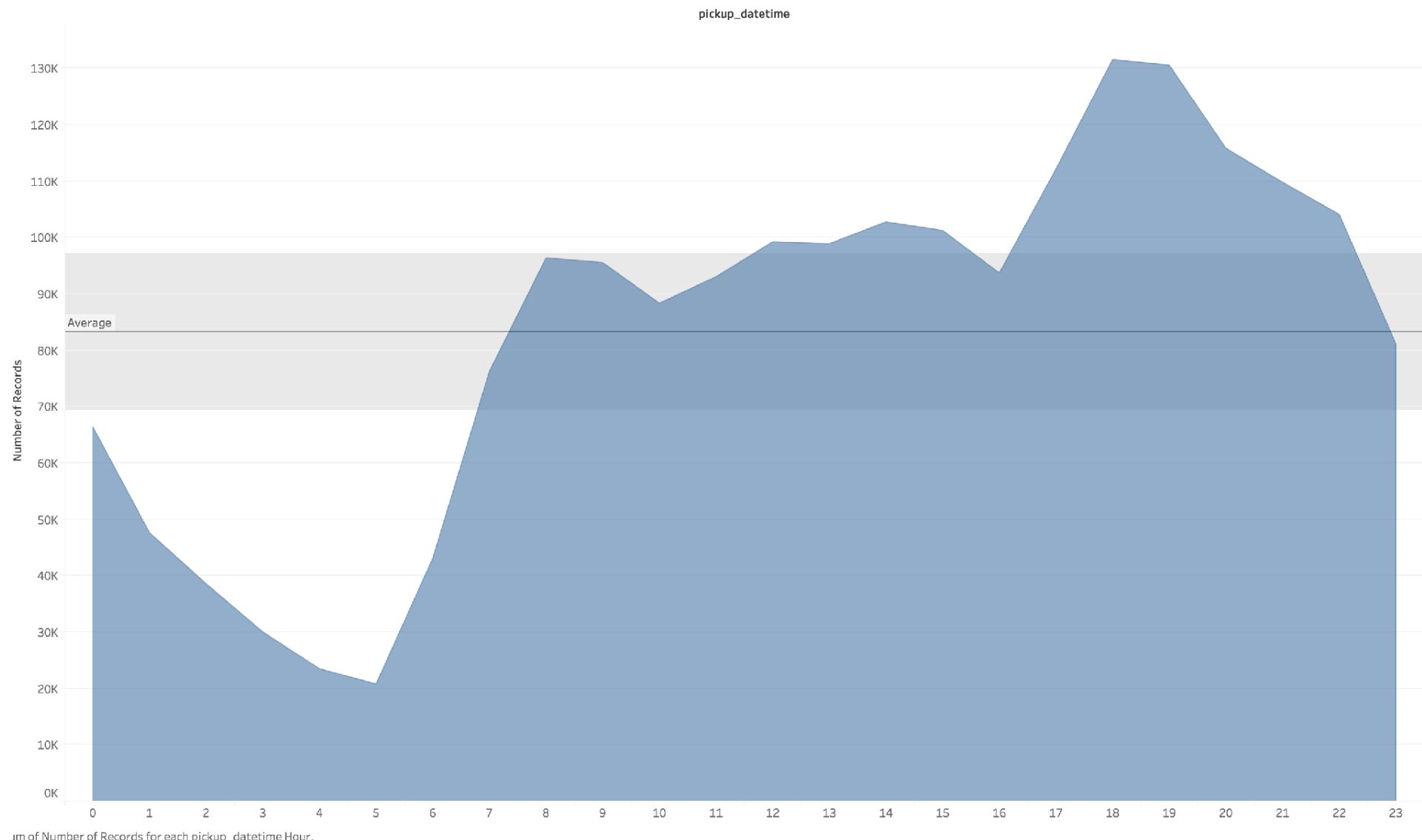
pickup\_datetime



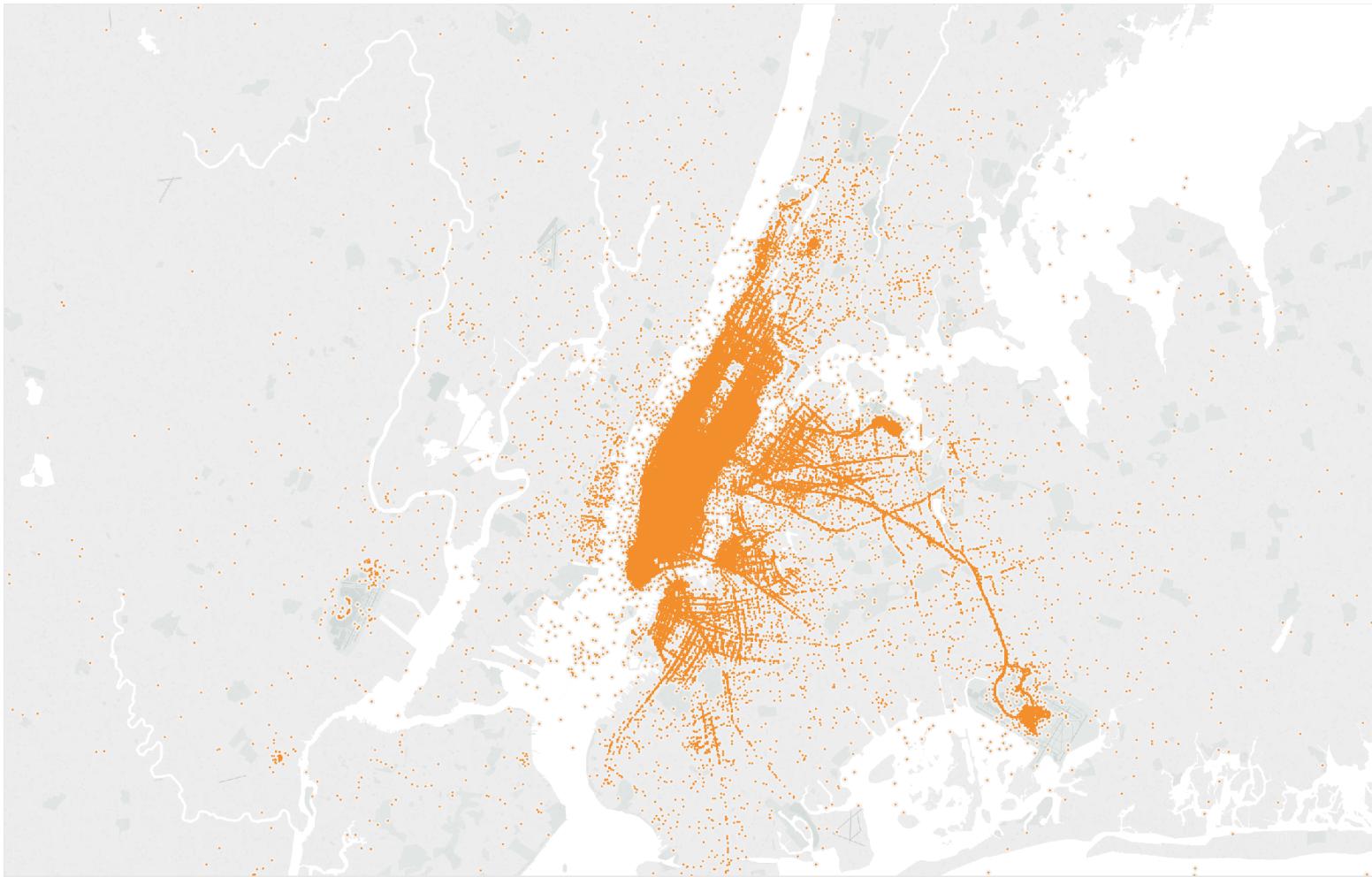
Sum of Number of Records for each pickup\_datetime Day.

## Distribution of number of trips per hour

trips per day

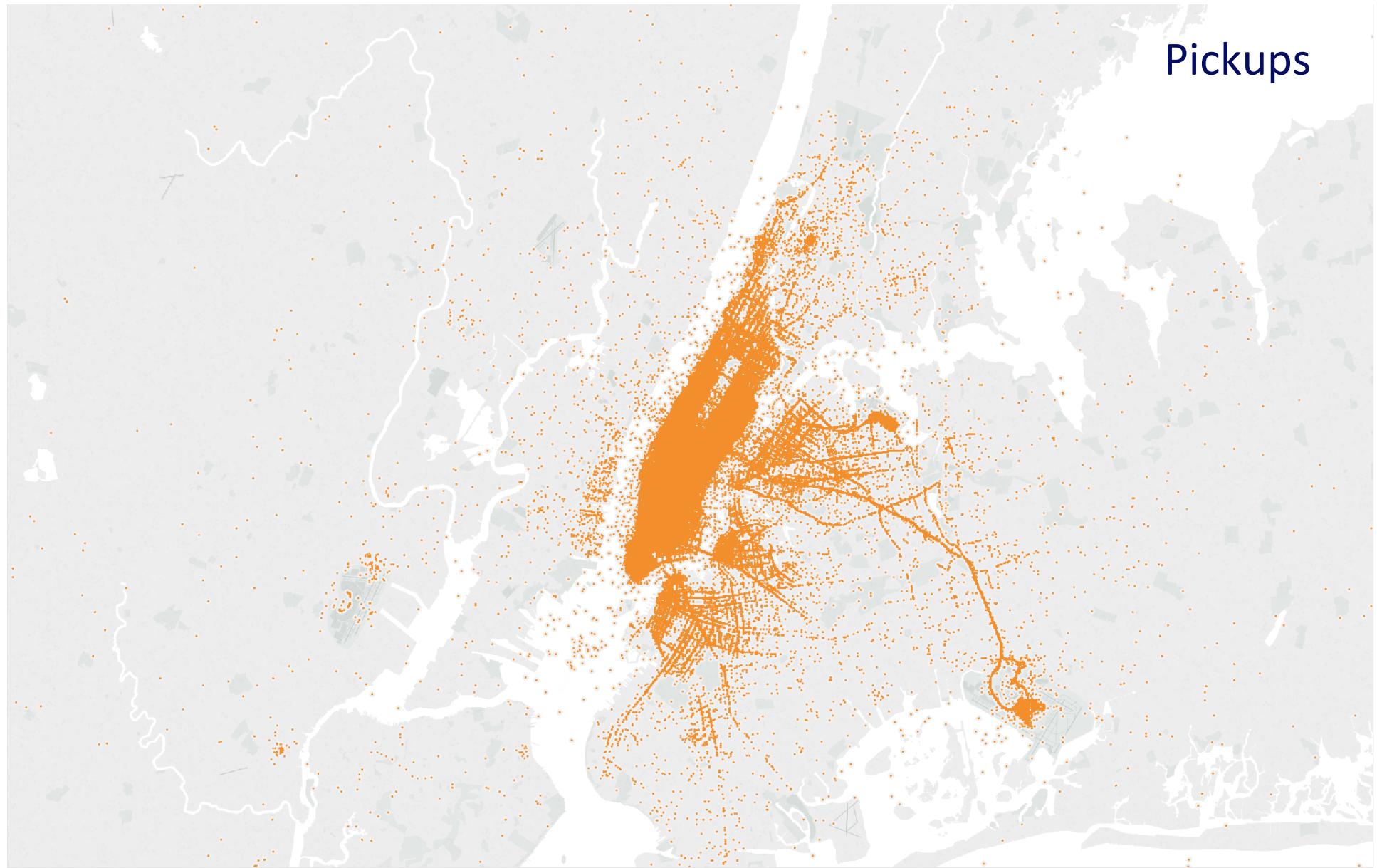


# Data for 20 days: Pickups

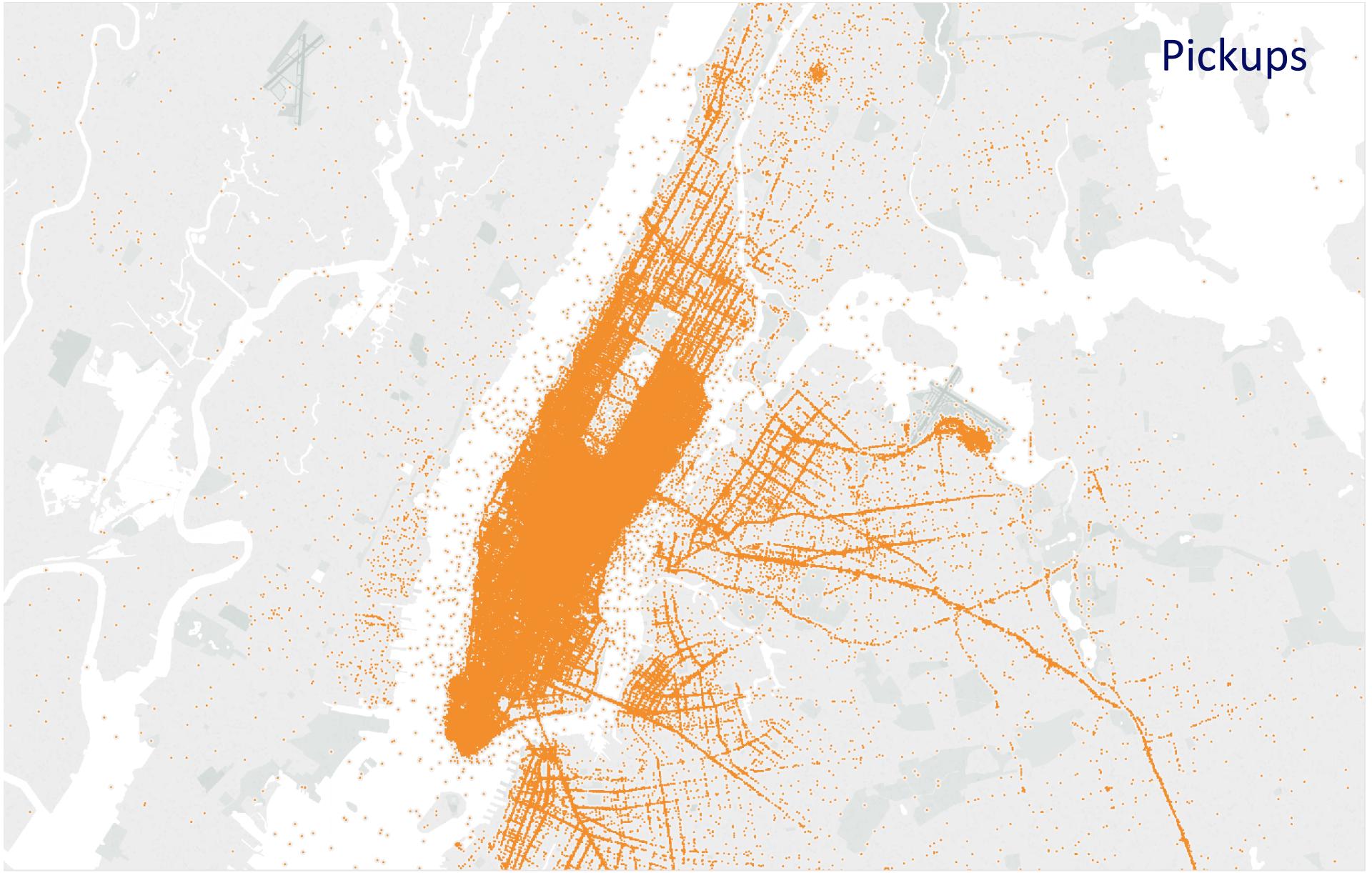


Map based on pickup\_longitude and pickup\_latitude. The data is filtered on pickup\_datetime, which ranges from 01/01/2013 00:00:00 to 20/01/2013 23:59:27.

# Pickups

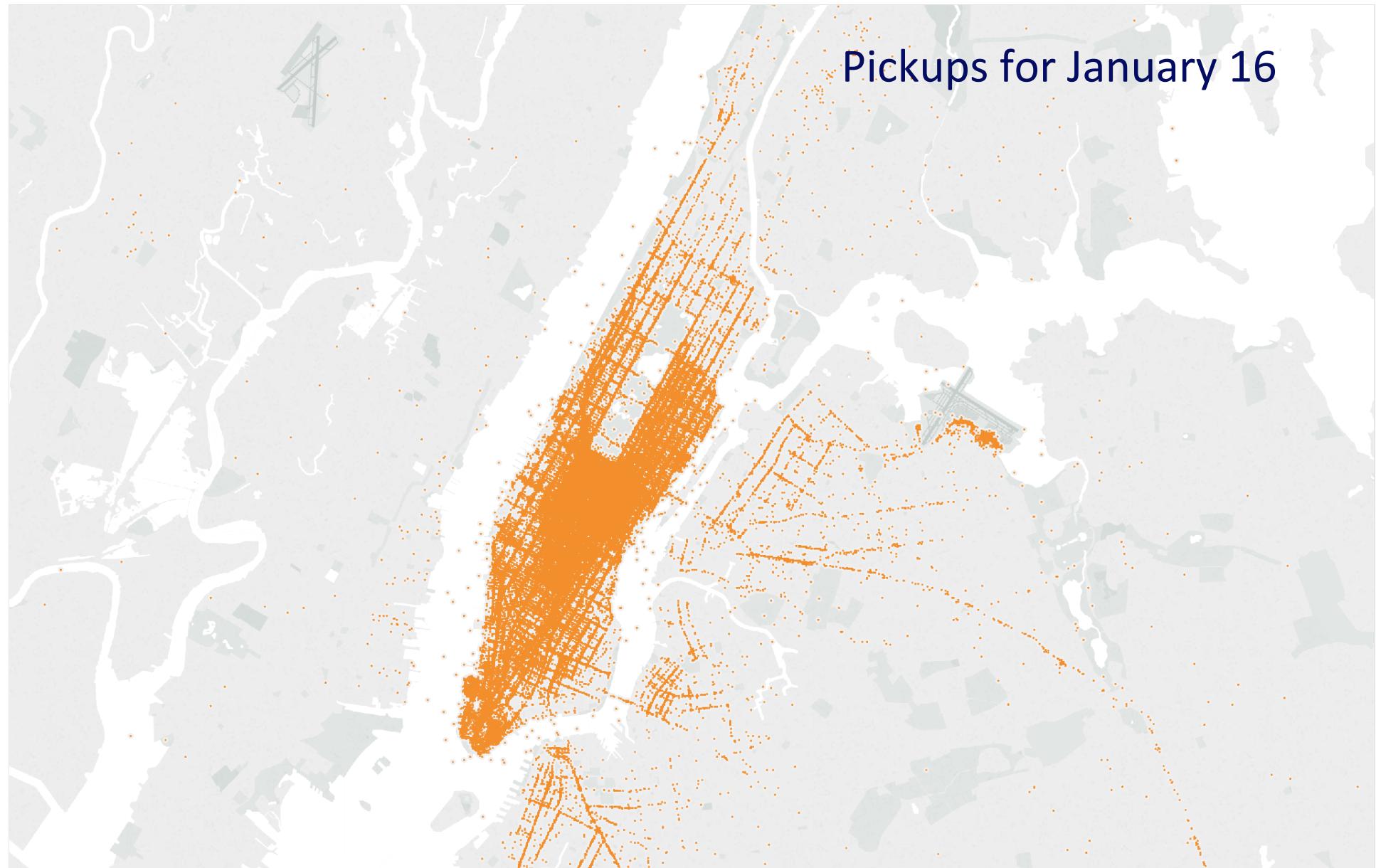


Map based on pickup\_longitude and pickup\_latitude. The data is filtered on pickup\_datetime, which ranges from 01/01/2013 00:00:00 to 20/01/2013 23:59:27.



# Pickups

Map based on pickup\_longitude and pickup\_latitude. The data is filtered on pickup\_datetime, which ranges from 01/01/2013 00:00:00 to 20/01/2013 23:59:27.

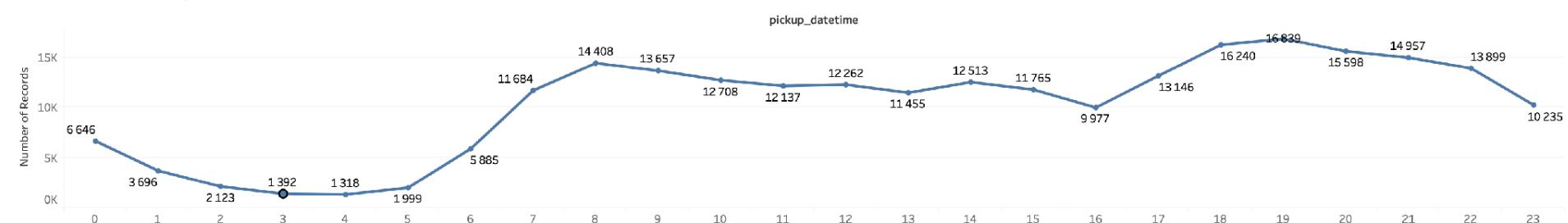


Map based on pickup\_longitude and pickup\_latitude. The data is filtered on pickup\_datetime, which ranges from 16/01/2013 00:00:00 to 16/01/2013 23:59:27.

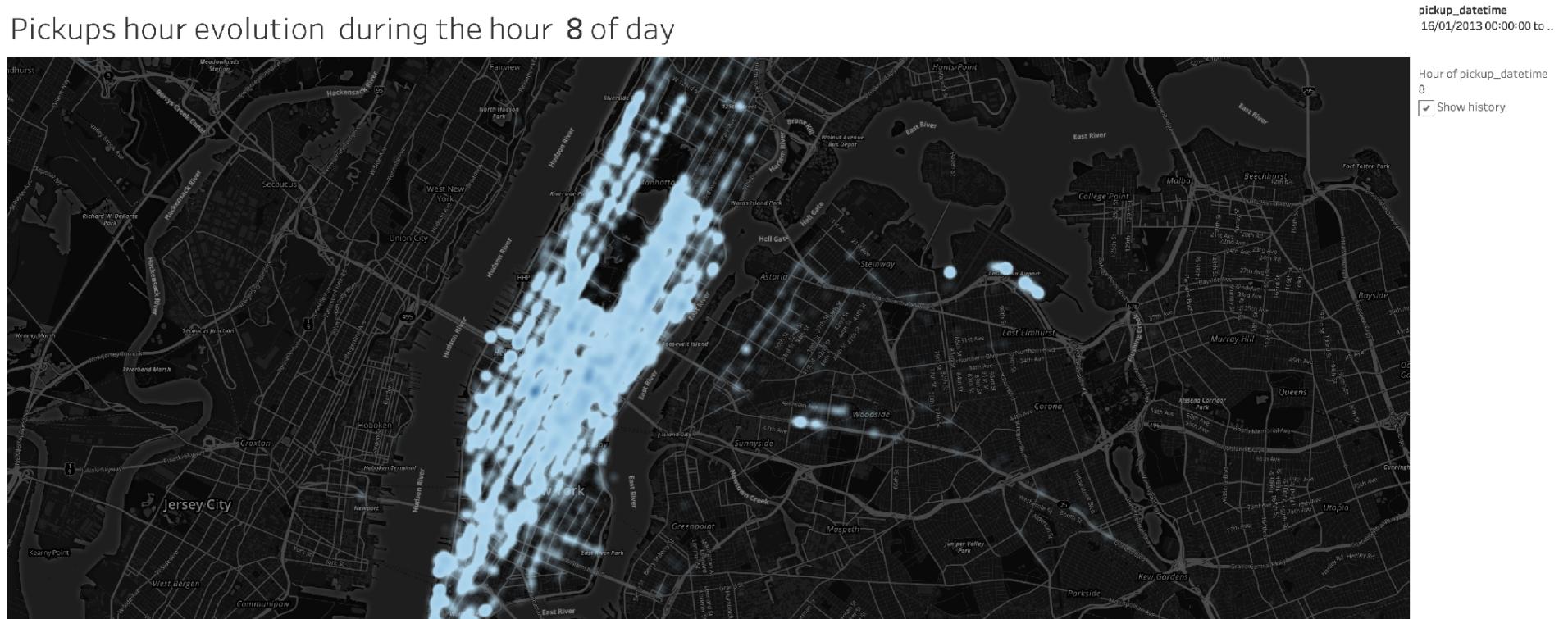
## Pickups hour evolution during the hour 3 of day



## Hour Evolution a day



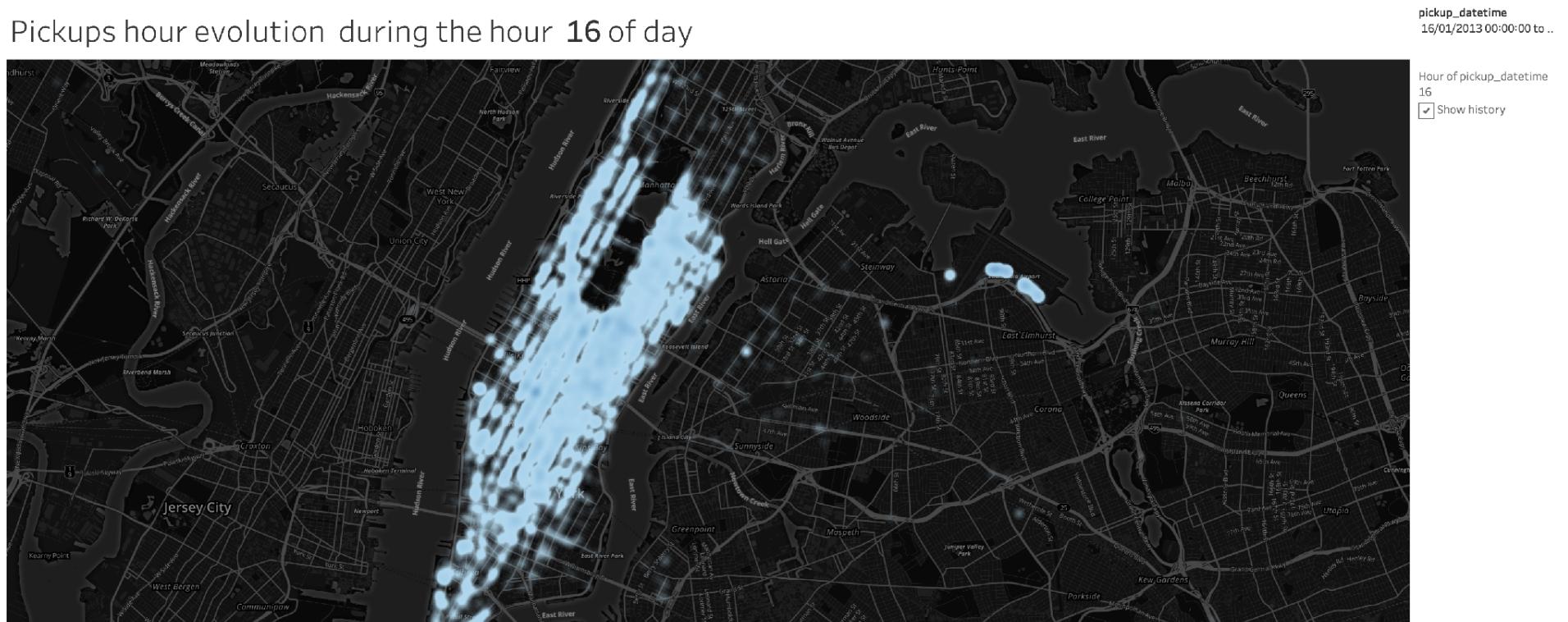
## Pickups hour evolution during the hour 8 of day



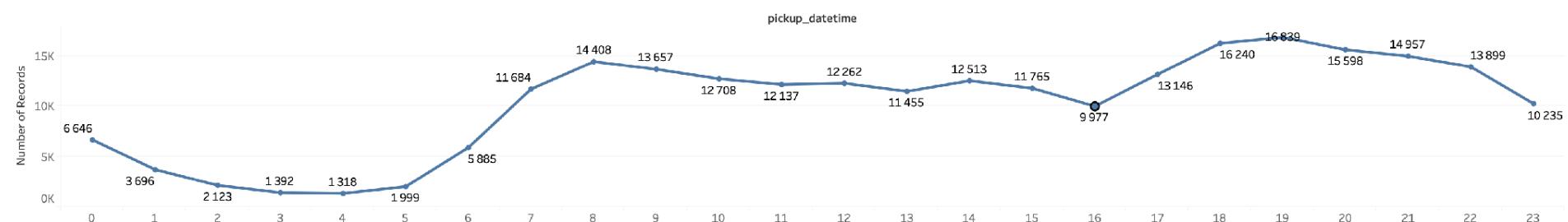
## Hour Evolution a day



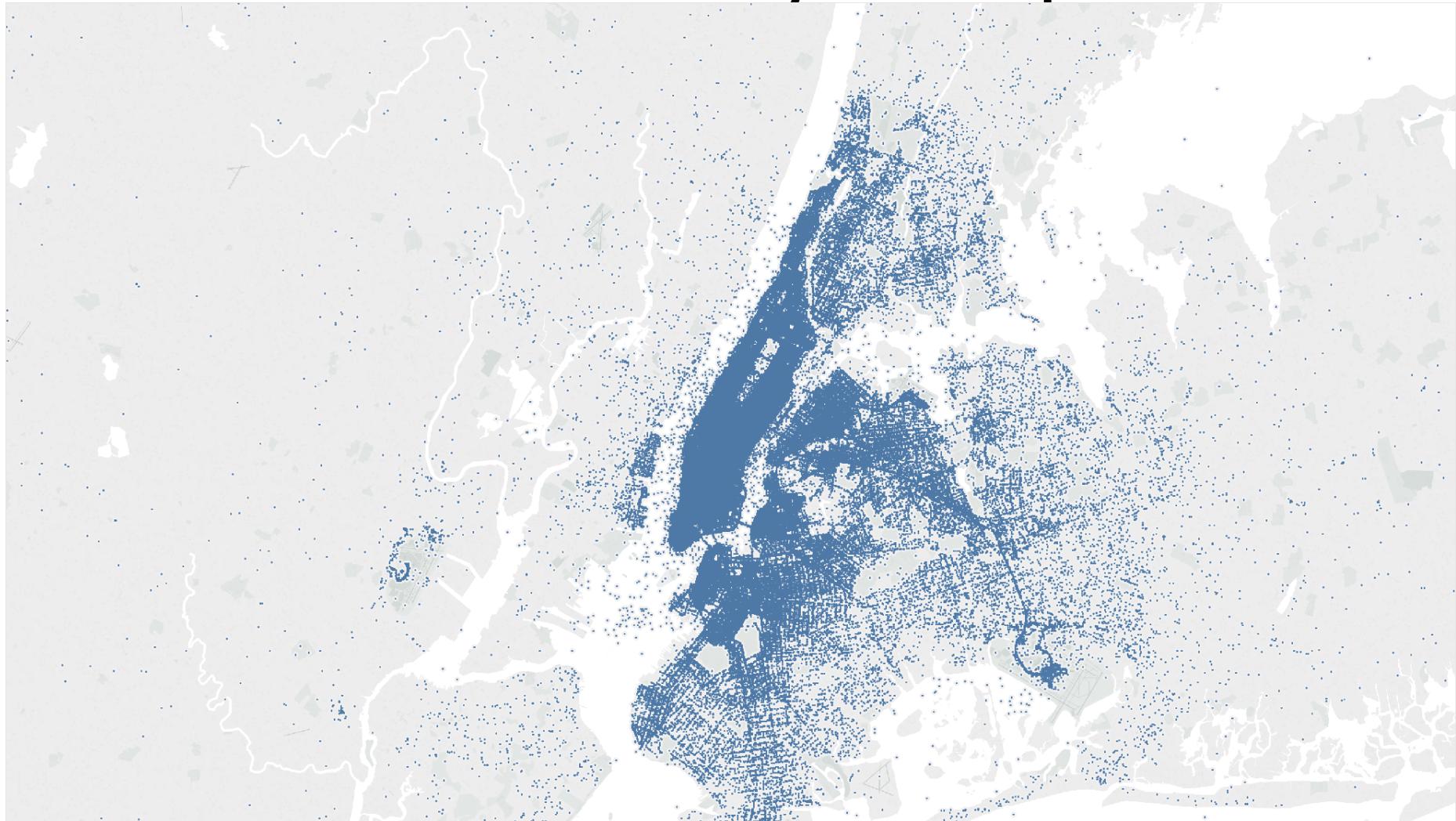
## Pickups hour evolution during the hour 16 of day



Hour Evolution a day



# Data for 20 days: Drop-off



Map based on dropoff\_longitude and dropoff\_latitude.

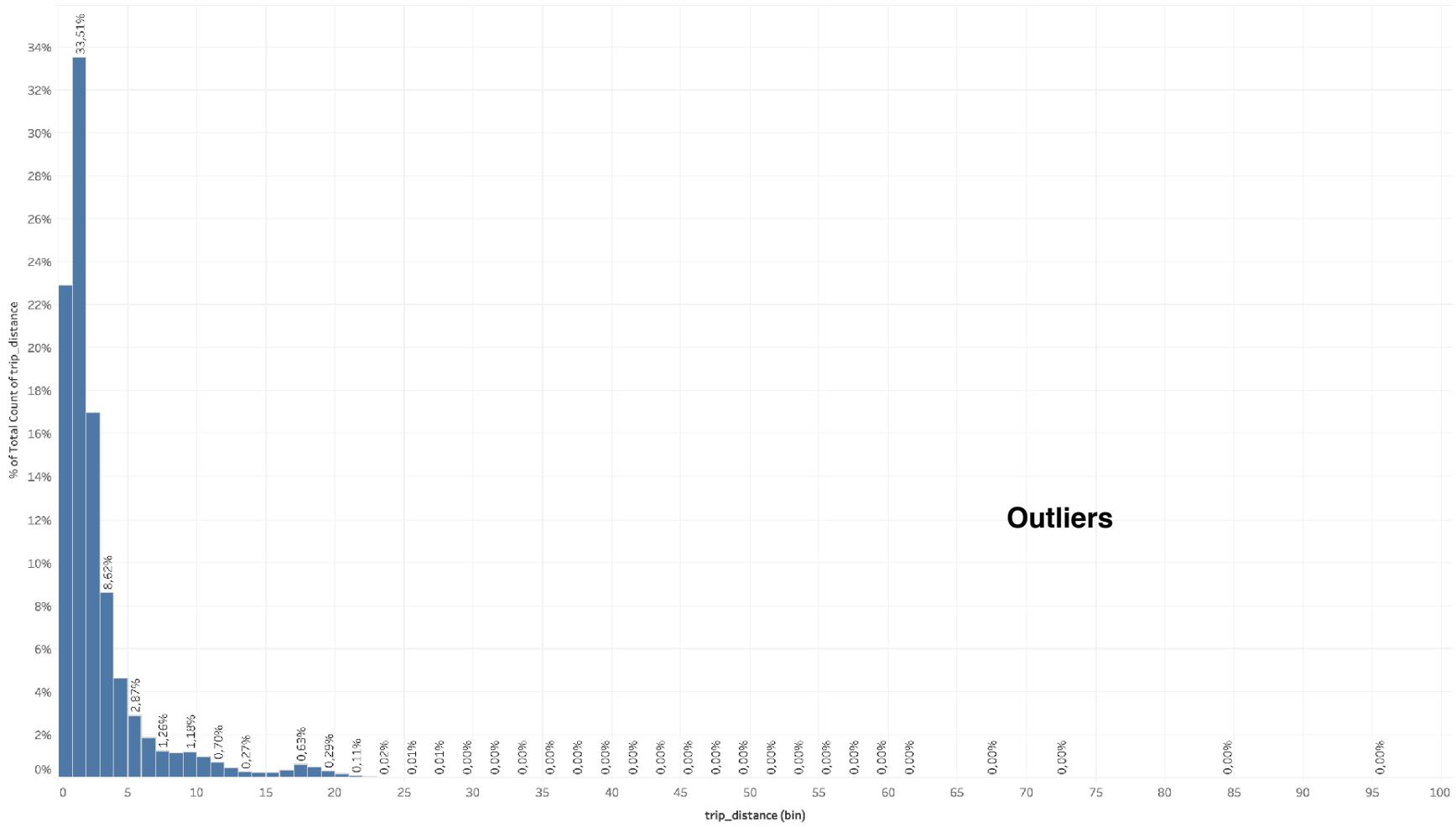
Pickups

New York City  
2009-2015

Drop offs



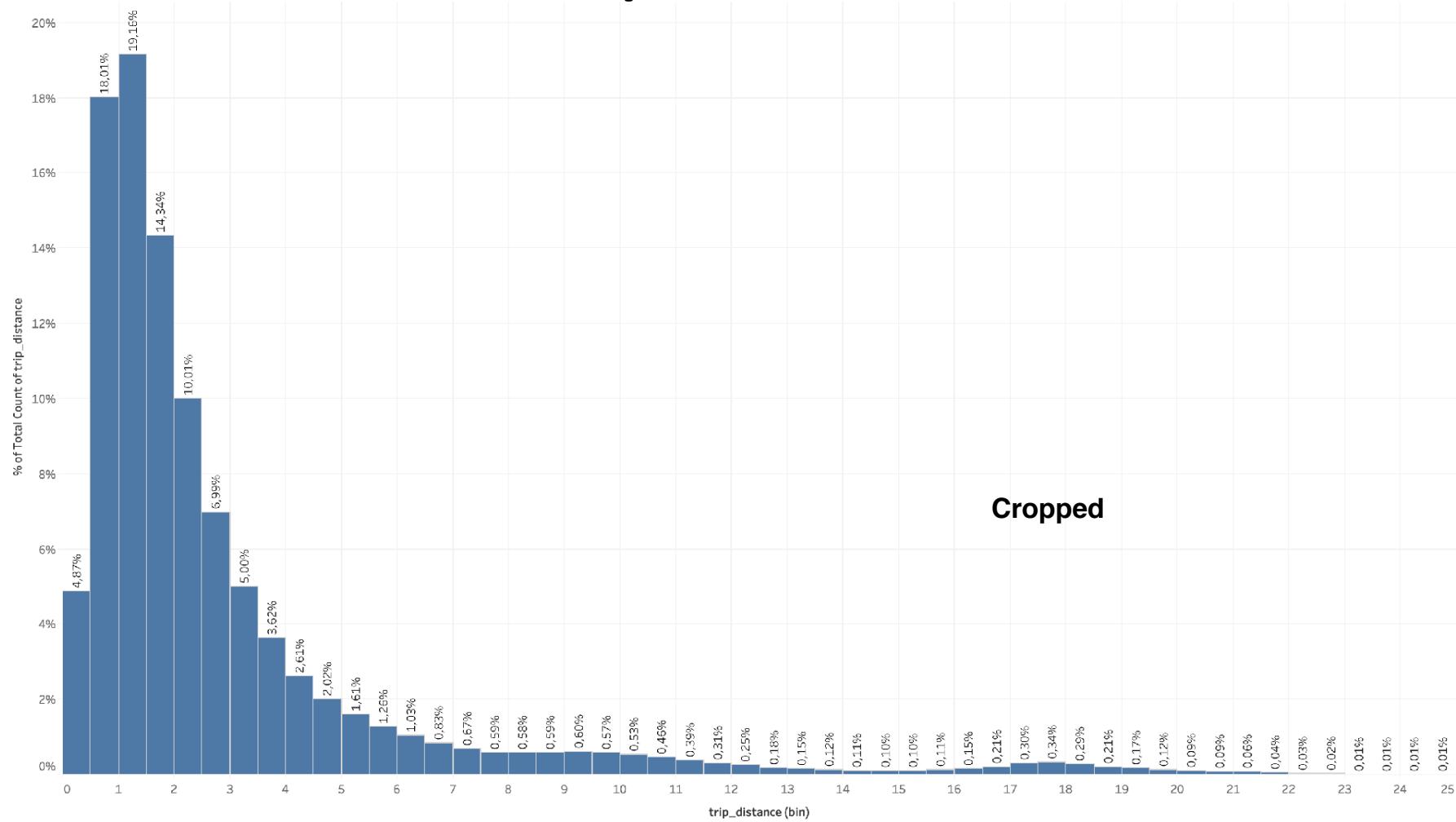
# Trip Distance



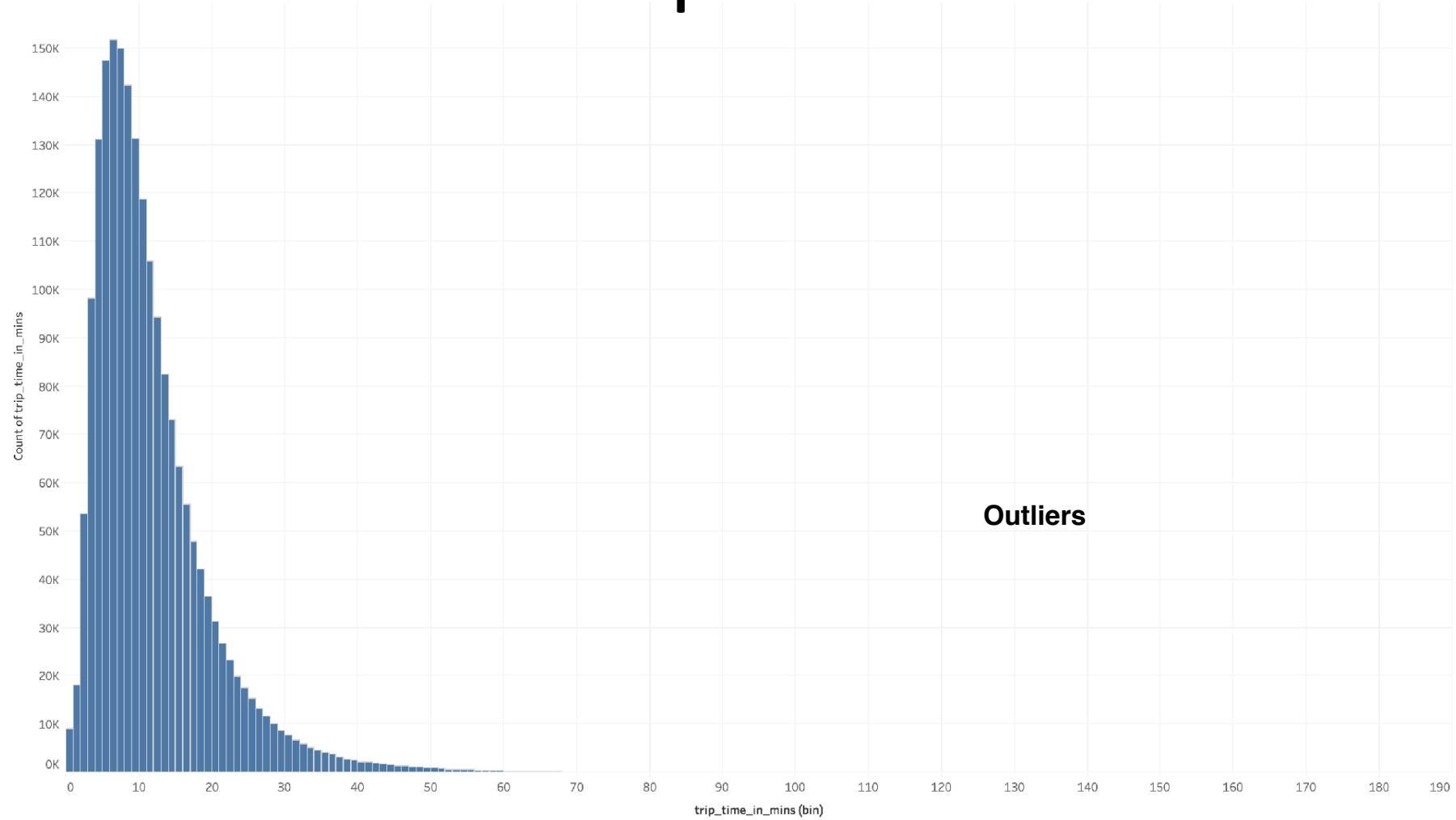
Outliers

The trend of % of Total Count of trip\_distance for trip\_distance (bin).

# Trip Distance

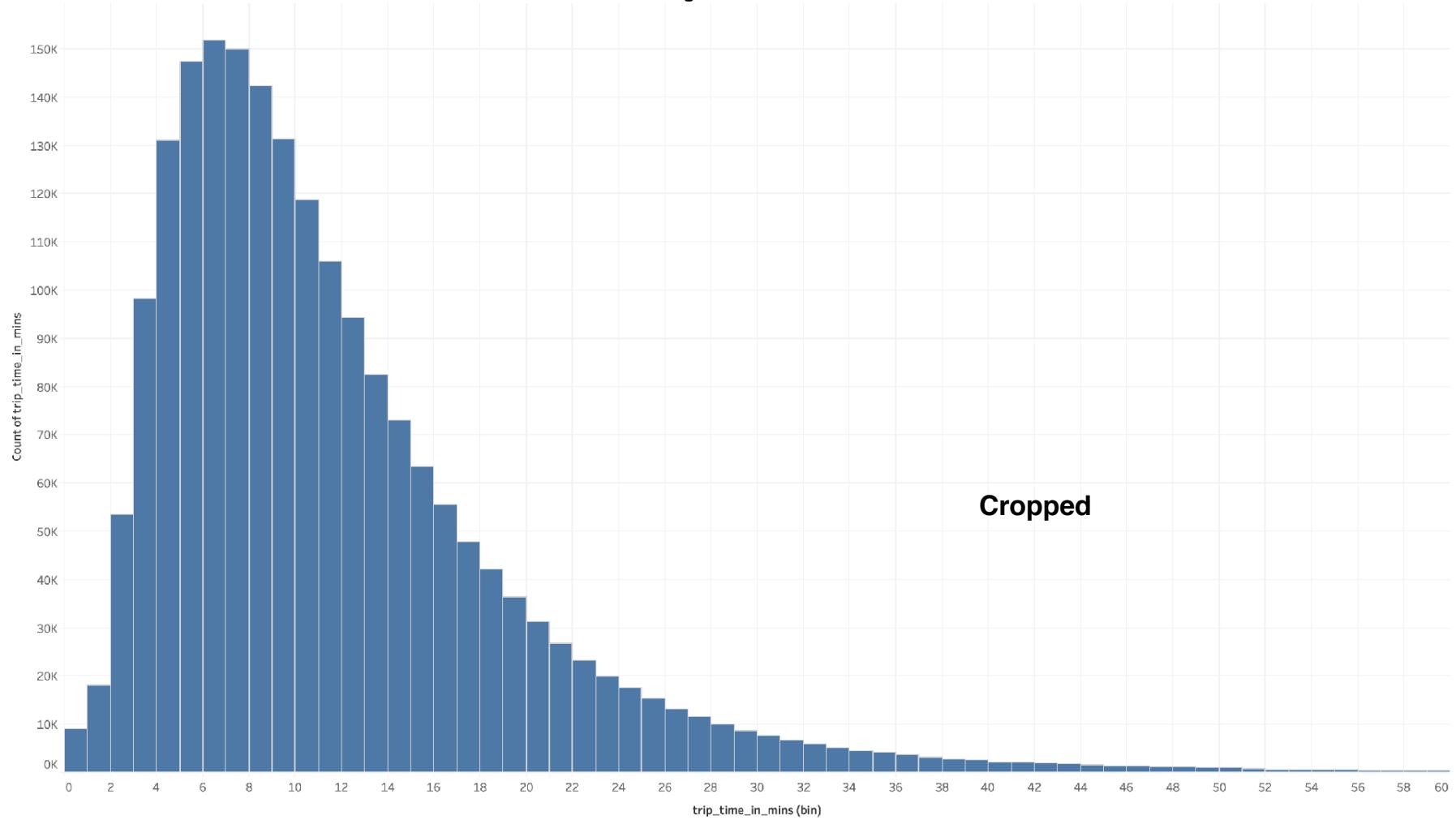


# Trip Time



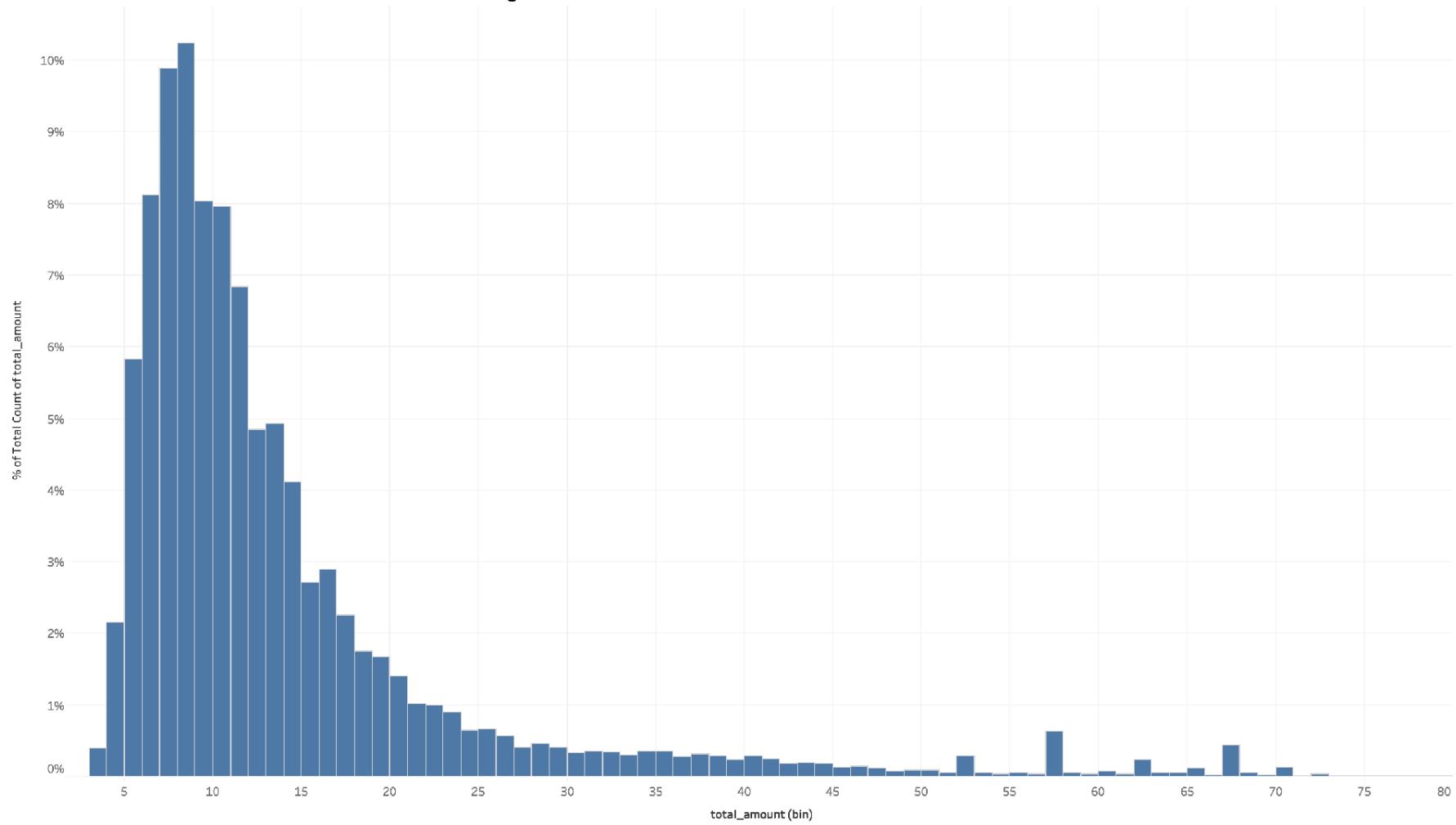
The trend of count of trip\_time\_in\_mins for trip\_time\_in\_mins (bin).

# Trip Time



The trend of count of trip\_time\_in\_mins for trip\_time\_in\_mins (bin).

# Trip Total Amount



The trend of % of Total Count of total\_amount for total\_amount (bin).

# Full Assignment

- Full assignment text and supplemental code and scripts:

<https://github.com/smduarte/ps2021>