

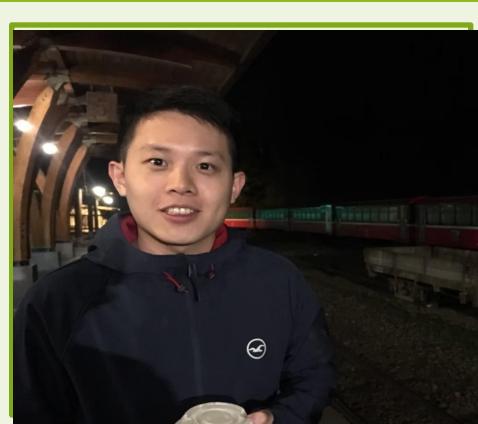


Exploring the Usage of Divvy Bikes

PRESENTED BY:

McKenzie Campbell, Louise Wang, Howard Lin

Team



Howard Lin



Louise Wang



McKenzie Campbell

Business Cases

Business Use Case	Data Processing	Data Modeling	Visualization	Recommendations
-------------------	-----------------	---------------	---------------	-----------------

- **Identifying** factors that affect bike usage on a localized level
- **Optimizing** the allocation of bikes numbers to each bike station
- **Building** new Divvy stations in areas with potential for growth



Executive Summary

- The purpose of this system is to support decision making process for Divvy Bike company including resource and infrastructure planning.
- By extracting, organizing, analyzing, visualizing data using a variety of platforms, we aim to identify the factors contributing to the usage and trends of Divvy Bike in Chicago.
- System functionality includes ETL process from raw data collectors, relational database for storage and access, BI reporting, connections to visualization tools.

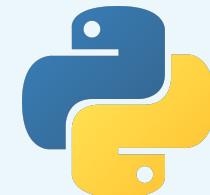
Data Profile

Time Horizon:
April 2020 - October 2020

	Description	Data Size	Shape	Key Variables	Matching Fields
Divvy Trips	Data broken down for each trip taken	814.4 MB	6,870,685 rows, 10 columns	start_station, start_datetime, end_station, end_datetime	Date, hour
Divvy Stations	Geographical information on each Divvy Bike Station	66 KB	669 row, 8 columns	station_geometry, station_total_docks	Zip code, region
Weather	Weather taken from Midway broken down by hour	9.6 MB	23,443 rows, 20 columns	temperature, weather type	Hour
Traffic	Traffic for each neighborhood in Chicago	799.1 MB	1,048,575 rows, 17 columns	region, ave_speed	Hour, region
Covid Cases	Cases by zip code in Chicago	26.6 MB	183,317 rows, 48 columns	zip codes, case measures	Date, zip code
Region	Geographical information for each neighborhood	5 KB	29 rows 3 columns	geometry	Region

Tools

Data Collection



Data preparation



WORKS WITH
MySQL™



Refine^{OPEN}

Data Storage



Google Cloud

WORKS WITH
MySQL™

NoSQL



mongoDB

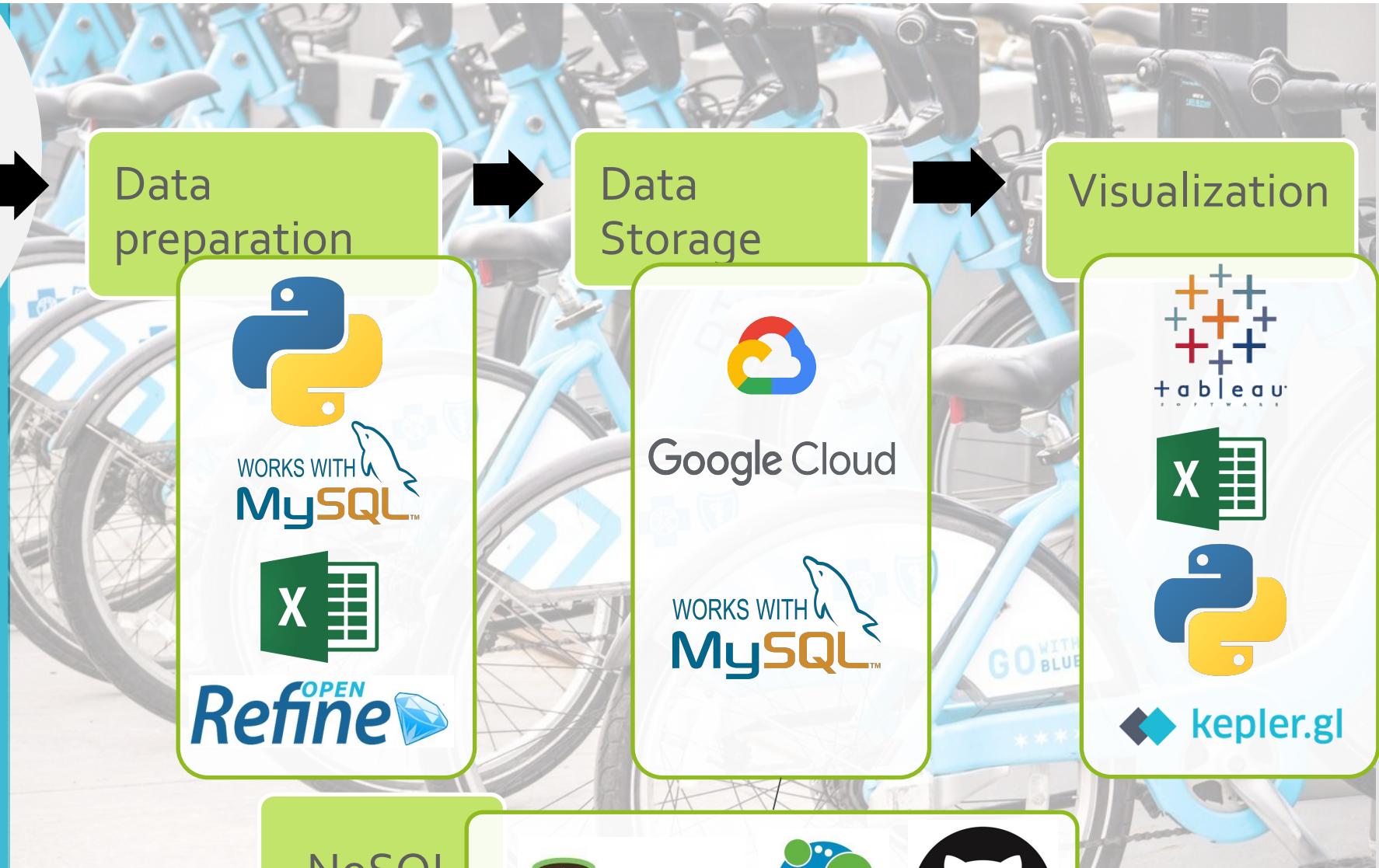


GitHub

Visualization



kepler.gl



Data Preparation [Trip]



Jupyter Notebook

```
##cleaned.trips.Q2.2019.csv##  
  
read_file = pd.read_csv ('/Users/KenzieKay/Desktop/Engineering Platforms MSCA 31012/1  
read_file.to_csv ('/Users/KenzieKay/Desktop/Engineering Platforms MSCA 31012/Trip Dat  
  
data10 = pd.read_csv('/Users/KenzieKay/Desktop/Engineering Platforms MSCA 31012/Trip  
  
data10 = data10.rename(columns={'01 - Rental Details Rental ID':'ride_id'})  
data10.drop('01 - Rental Details Bike ID', inplace=True, axis=1)  
data10.drop('01 - Rental Details Duration In Seconds Uncapped', inplace=True, axis=1)  
data10 = data10.rename(columns={'03 - Rental Start Station ID':'start_station_id'})  
data10 = data10.rename(columns={'03 - Rental Start Station Name': 'start_station_name'})  
data10 = data10.rename(columns={'02 - Rental End Station ID':'end_station_id'})  
data10 = data10.rename(columns={'02 - Rental End Station Name': 'end_station_name'})  
data10 = data10.rename(columns={'User Type':'member_casual'})  
data10.drop('Member Gender', inplace=True, axis=1)  
data10.drop('05 - Member Details Member Birthday Year', inplace=True, axis=1)  
  
data10['start_station_name'].replace('', np.nan, inplace=True)  
data10['start_station_id'].replace('', np.nan, inplace=True)  
data10['end_station_name'].replace('', np.nan, inplace=True)  
data10['end_station_id'].replace('', np.nan, inplace=True)  
  
data10.dropna(subset=['start_station_name'], inplace=True)  
data10.dropna(subset=['start_station_id'], inplace=True)  
data10.dropna(subset=['end_station_name'], inplace=True)  
data10.dropna(subset=['end_station_id'], inplace=True)  
  
new10a = data10["01 - Rental Details Local Start Time"].str.split(" ", n = 1, expand = True)  
data10[ "started_date" ]= new10a[0]  
data10[ "started_time" ]= new10a[1]  
data10.drop(columns =["01 - Rental Details Local Start Time"], inplace = True)  
  
new10b = data10["01 - Rental Details Local End Time"].str.split(" ", n = 1, expand = True)  
data10[ "ended_date" ]= new10b[0]  
data10[ "ended_time" ]= new10b[1]  
data10.drop(columns =["01 - Rental Details Local End Time"], inplace = True)
```

Tools and Methods Utilized

Jupyter Notebook

- Drop columns that were not across files
- Dropped rows that did not have a start or stop station
- Rename columns to match
- Read to CSV for files that were originally txt files.
- **12 original files**
- **Ended up with over 6 million rows**



Data Preparation [Station]



Jupyter Notebook

```
!pip install geopandas
...
import pandas as pd
import geopandas as gpd
stationpoints = pd.read_csv('/Users/howardlin/Desktop/DEPA data/Divvy_Bicycle_Stations.csv', skiprows=0, encoding='utf-8')
...
stationpoints=stationpoints.drop(columns='Location')
gdf_station = gpd.GeoDataFrame(stationpoints,
                                geometry=gpd.points_from_xy(stationpoints.Longitude,stationpoints.Latitude),
                                crs="EPSG:4326")
...
gdf_zip = gpd.read_file('/Users/howardlin/Desktop/DEPA data/ZIP_Codes.geojson')
...
sjoined_station = gpd.sjoin(gdf_station, gdf_zip, op="within")
```

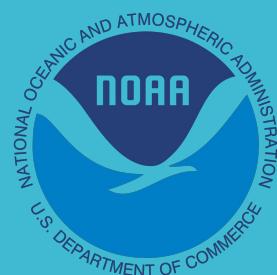
Output

	A	B	C	D	E	F	G	H	I	J	K
1	ID	Station Name	Total Docks	Docks in Ser	Status	geometry	zip	REGION_ID	REGION		
2	715	Major Taylor Trail & 124th St	15	15	In Service	POINT (-87.63711333299999 41.668696353)	60628	28	Riverdale-Hegewisch		
3	695	Dauphin Ave & 103rd St	15	15	In Service	POINT (-87.60766267841.707250368)	60628	26	Washington Hts-Roseland-Pullman		
4	698	Cottage Grove Ave & 111th Pl	15	15	In Service	POINT (-87.61000156441.691709758)	60628	26	Washington Hts-Roseland-Pullman		
5	699	Vernon Ave & 107th St	15	14	In Service	POINT (-87.61225998399999 41.700121379)	60628	26	Washington Hts-Roseland-Pullman		
6	700	Halsted St & 96th St	15	15	In Service	POINT (-87.64301955741.719712184)	60628	26	Washington Hts-Roseland-Pullman		
7	703	Greenwood Ave & 97th St	15	15	In Service	POINT (-87.59703576641.718703155)	60628	26	Washington Hts-Roseland-Pullman		
8	1441	S Wentworth Ave & W 111th St	5	5	In Service	POINT (-87.628051799999941.6926176000001)	60628	26	Washington Hts-Roseland-Pullman		
9	775	Halsted St & 108th St	15	15	In Service	POINT (-87.612121070000141.7051207801)	60628	26	Washington Hts-Roseland-Pullman		

Tools and Methods Utilized

Jupyter Notebook (packages pandas,geopandas): Given the geojson file for area boundaries, we use Spatial Join in Geopandas to match each Divvy Bike station into **Zip Code** and **Neighborhood region**.

Data Preparation [Weather]



Jupyter Notebook

```
import pandas as pd

data = pd.read_csv('/Users/howardlin/Documents/Class documents/Data Engineering /team4 project/data/Local_climatological_data_chicago_midway_daily.csv', skiprows=1)

data['DATE'] = pd.to_datetime(data['DATE'])

data['DAY_OF_WEEK'] = data.DATE.dt.dayofweek

data.rename(columns={'DATE': 'DATETIME'}, inplace=True)

data['DATE'] = data.DATETIME.dt.date

data['HOUR'] = data.DATETIME.dt.hour

data["TIME_CHECK"] = data["DATE"].astype(str) + " " + data["HOUR"].astype(str)

data.REPORT_TYPE.value_counts()

data[~data.DATE.str.contains('2019-00')] = None
```

Output

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U			
1	DATETIME	REPORT_TYF	SOURCE	HourlyAltime	HourlyDewPt	HourlyDryBul	HourlyPrecip	HourlyPreser	HourlyPressu	HourlyRelati	HourlySeale	HourlySkyCor	HourlyStatio	HourlyVisibil	HourlyWetBt	HourlyWindC	HourlyWindG	HourlyWinds	DAY_OF_WEEK	DATE			
2	2019/1/1 0:53	FM-15		7	29.98	33	36	0	0	0	89	30	BKN:07 14 0	29.32	10	35	320	Null	10	1	2019/1/1		
3	2019/1/1 1:53	FM-15		7	30.03	32	36	0	0	0	86	30.06	BKN:07 14 0	29.37	10	34	330	Null	13	1	2019/1/1		
4	2019/1/1 2:53	FM-15		7	30.08	32	35	0	-0.12	3	89	30.1	BKN:07 12 B	29.41	8	34	340	Null	9	1	2019/1/1		
5	2019/1/1 3:53	FM-15		7	30.11	31	34	0.005	Drizzle	0	0	89	30.14	OVC:08 13	29.44	7	33	330	Null	10	1	2019/1/1	
6	2019/1/1 4:53	FM-15		7	30.13	31	33	0.005	Drizzle	Mist	0	0	92	30.16	BKN:07 7 B	29.47	5	32	360	Null	7	1	2019/1/1
7	2019/1/1 5:53	FM-15		7	30.17	29	32	0.005	Drizzle	Mist	-0.09	3	88	30.2	SCT:04 8 BK	29.5	5	31	20	Null	11	1	2019/1/1
8	2019/1/1 6:53	FM-15		7	30.22	27	31	0.005	Drizzle	Mist	0	0	85	30.25	OVC:08 14	29.55	6	30	10	20	15	1	2019/1/1
9	2019/1/1 7:53	FM-15		7	30.24	26	31	0	0	0	82	30.27	OVC:08 16	29.57	10	29	10	Null	11	1	2019/1/1		
10	2019/1/1 8:53	FM-15		7	30.27	27	31	0	-0.09	1	85	30.3	OVC:08 14	29.6	10	30	350	Null	9	1	2019/1/1		
11	2019/1/1 9:53	FM-15		7	30.29	27	31	0.005	Drizzle	0	0	85	30.32	OVC:08 13	29.62	7	30	340	Null	9	1	2019/1/1	
12	2019/1/1 10:53	FM-15		7	30.29	28	31	0.005	Drizzle	Mist	0	0	89	30.32	BKN:07 10 0	29.62	6	30	330	Null	9	1	2019/1/1
13	2019/1/1 11:53	FM-15		7	30.28	28	31	0.005	Drizzle	-0.02	0	89	30.32	OVC:08 12	29.61	7	30	360	Null	11	1	2019/1/1	
14	2019/1/1 12:53	FM-15		7	30.28	26	30	0.005	Drizzle	0	0	85	30.32	OVC:08 11	29.61	8	29	340	Null	10	1	2019/1/1	
15	2019/1/1 13:53	FM-15		7	30.28	25	29	0.005	Drizzle	0	0	85	30.32	OVC:08 10	29.61	7	28	350	Null	8	1	2019/1/1	
16	2019/1/1 14:53	FM-15		7	30.29	25	29	0.005	Drizzle	0	3	85	30.32	OVC:08 10	29.62	8	28	20	Null	8	1	2019/1/1	
17	2019/1/1 15:53	FM-15		7	30.31	25	29	0.005	Drizzle	0	0	85	30.34	OVC:08 11	29.64	8	28	340	Null	8	1	2019/1/1	
18	2019/1/1 16:53	FM-15		7	30.32	24	28	0.005	Drizzle	0	0	85	30.35	BKN:07 12 0	29.65	9	27	340	Null	6	1	2019/1/1	
19	2019/1/1 17:53	FM-15		7	30.27	24	28	0.005	Drizzle	-0.02	1	85	30.36	CFT:04 12 BK	29.65	9	27	330	Null	7	1	2019/1/1	

Tools and Methods Utilized

Jupyter Notebook (packages pandas, datetime) to get columns, rows and format needed.

OpenRefine to translate technical weather code in the data.

Business Use Case

Data Processing

Data Modeling

Visualization

Recommendations

OpenRefine

The screenshot shows the OpenRefine interface with a dataset titled "Local_climatological_data_chicago_midway.csv". The interface includes a header row with columns: STATION, DATE, REPORT_TYPE, SOURCE, HourlyAltime, HourlyDewPt, HourlyDryBul, HourlyPrecip, HourlyPreser, HourlyPressu, HourlyRelati, HourlySeale, HourlySkyCor, HourlyStatio, HourlyVisibil, HourlyWetBt, HourlyWindC, HourlyWindG, HourlyWinds, DAY_OF_WEEK, and DATE. The data table below contains approximately 23,443 rows of weather data. The interface also shows various filters and search tools at the top.

Data Preparation [Traffic]



Jupyter Notebook

```

import pandas as pd
import numpy as np

ita = pd.read_csv('/Users/howardlin/Downloads/Chicago_Traffic_Tracker_-_Historical_Congestion_Estimates_by_Regions.csv')

ita['TIME1'] = pd.to_datetime(data['TIME'])

ita['DATE'] = data.TIME1.dt.date

ataselect = data[(data["TIME1"] > '2019-01-01 00:00:00') & (data["TIME1"] <= '2020-11-15 23:59:59')]

dataselect = dataselect.drop(columns=['WEST', 'EAST', 'SOUTH', 'NORTH'])

dataselect = dataselect.sort_values(by=['TIME1', 'REGION_ID'])

dataselect['TIME'] = pd.to_datetime(dataselect['TIME'])

dataselect = dataselect.drop(columns=['TIME1'])

dataselect.to_csv('/Users/howardlin/Desktop/DEPA data/traffic.csv',
                 index=False)

all_together = (dataselect.groupby(['DATE', 'MONTH', 'DAY_OF_WEEK', 'HOUR', 'REGION_ID', 'REGION', 'NW_LOCATION', 'SE_L'])
               .agg({'SPEED': [np.mean, np.min, np.max]})  

               .rename(columns={'mean': 'avg_speed', 'amin': 'min_speed', 'amax': 'max_speed'}))

all_together = all_together.reset_index()

```

OpenRefine

469970 rows

	date	hour	region	avg_speed	min_speed	max_speed
1	2019-01-01	0	Rogers Park - West Ridge	28.55	25.98	30
2	2019-01-01	1	Far North West	25.04	23.73	27.95
3	2019-01-01	1	North Park-Albany-Lincoln Sq	26.42	24.55	28.64
4	2019-01-01	1	Edge Water-Uptown	24.44	21.62	26.26
5	2019-01-01	1	Dunning-Portage-Belmont-Cragin	26.42	23.23	28.88
6	2019-01-01	1	Irving Park-Avondale-North Ctr	26.73	23.91	27.27
7	2019-01-01	1	Humboldt Lagoon Square	24.94	23.98	26.82
8	2019-01-01	1	UIC/Loyola Lakeview	25.96	21.34	26.95
9	2019-01-01	1	Austin	21.46	6.0	26.92
10	2019-01-01	1	West Garfield Park-Sixth Ward	26.68	25.91	27.95
11	2019-01-01	1	West Town-Avondale	21.1	13.23	26.86
12	2019-01-01	1	West North Side	26.48	17.32	25.18
13	2019-01-01	1	Chicago Loop	16.84	16.0	21.82
14	2019-01-01	1	Lakeview	26.73	23.64	30.37
15	2019-01-01	1	Bricktown-Millenium-Lower Wabash	26.59	23.92	26.94
16	2019-01-01	1	Inner South-Douglas	23.99	17.31	26.82
17	2019-01-01	1	Milwaukee-Garfield-Pulaski	21.7	19.32	26.30
18	2019-01-01	1	South-West Side	23.98	24.88	26.88
19	2019-01-01	1	New City-Francisco-Vin Brinkwood	21.1	14.85	26.38
20	2019-01-01	1	Fisher Island-Brick-Ward-Park	21.75	22.01	27.28
21	2019-01-01	1	Hyde Park-Kingswood-Harold	21.97	23.23	26.99
22	2019-01-01	1	Wicker Park	25.56	6.0	30.99
23	2019-01-01	1	Austin-Greenwood-Chicago Park	26.06	24.02	27.27
24	2019-01-01	1	South Shore-E-Chicago-Jackson	28.41	24.96	32.98
25	2019-01-01	1	Bronx-Milwaukee-Morgan Park	16.02	6.0	35.43
26	2019-01-01	1	Washington-Haywood-Pulman	26.99	21.98	32.09
27	2019-01-01	1	South Chicago-Humboldt Park	24.75	20.0	32.18
28	2019-01-01	1	Riverside-Highland	24.15	6.0	35.18
29	2019-01-01	1	Columbus-Kingswood	16.19	6.0	33.9
30	2019-01-01	1	Lawndale-Lakeview	26.73	23.23	28.88
31	2019-01-01	1	Rogers Park - North Ridge	26.36	24.88	28.8
32	2019-01-01	1	Far North West	26.36	24.88	28.8
33	2019-01-01	1	North Park-Albany-Lincoln Sq	21.83	24.88	28.8

Output

	A	B	C	D	E	F	G	H	I
1	DATE	MONTH	DAY_OF_WEEK	HOUR	avg_speed	min_speed	max_speed	REGION_ID	REGION
2	2019/1/1		1	1	0	28.55	25.98	30	1 Rogers Park - West Ridge
3	2019/1/1		1	1	0	25.04	23.73	27.95	2 Far North West
4	2019/1/1		1	1	0	26.42	24.55	28.64	3 North Park-Albany-Lincoln Sq
5	2019/1/1		1	1	0	24.44	21.82	29.25	4 Edge Water-Uptown
6	2019/1/1		1	1	0	26.42	25.23	28.64	5 Dunning-Portage-Belmont-Cragin
7	2019/1/1		1	1	0	26.73	25.91	27.27	6 Irving Park-Avondale-North Ctr
8	2019/1/1		1	1	0	24.94	22.86	26.52	7 Humboldt Lagoon Square

Tools and Methods Utilized

Jupyter Notebook (packages pandas, numpy, datetime): We adjusted the time granularity from 10 min to hourly by Groupedby function.

OpenRefine filter out redundant row and column.

Data Preparation [Covid-19]



Jupyter Notebook

```
In [20]: covid.head()
```

```
Out[20]:
```

	date	zipcode	confirmed_cases	confirmed_cases_change	total_tested	total_tested_change	confirmed_cases_20
0	2020-11-27T00:00:00Z	60827	988		7	11878	91
1	2020-11-27T00:00:00Z	60707	2966		22	31312	294
2	2020-11-27T00:00:00Z	60661	444		3	7828	104
3	2020-11-27T00:00:00Z	60660	1484		16	28208	315
4	2020-11-27T00:00:00Z	60659	2181		18	21799	192

5 rows × 49 columns

```
In [21]: covid['date2'] = covid.date1.dt.date
```

```
In [15]: covid['day_of_week'] = covid.date1.dt.dayofweek
```

```
In [16]: covid['month'] = covid.date1.dt.month
```

[Microsoft Store]

Output

covid_record_id	date	late_update	month	day_of_week	weekend	zipcode	confirmed_cases_cumulative	confirmed_cases_increasing	tested_cases_cumulative	tested_cases_incre
1	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60827	988	7	11878	91
2	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60707	2966	22	31312	294
3	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60661	444	3	7828	104
4	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60660	1484	16	28208	315
5	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60659	2181	18	21799	192
6	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60657	2669	28	55367	699
7	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60656	1577	11	16076	104
8	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60655	1891	24	19074	163
9	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60654	1070	9	16474	227
10	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60653	1436	15	25050	181
11	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60652	3198	28	26009	217
12	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60651	4577	40	37243	297
13	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60649	1898	22	29708	223
14	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60647	4848	54	74772	747
15	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60646	1339	11	19142	135
16	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60645	2599	20	35393	317
17	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60644	2526	26	31269	133
18	2020-11-27	2020-12-05 03:50:55	11	4	Weekday	60643	2431	24	40567	317

Business Use Case

Data Processing

Data Modeling

Visualization

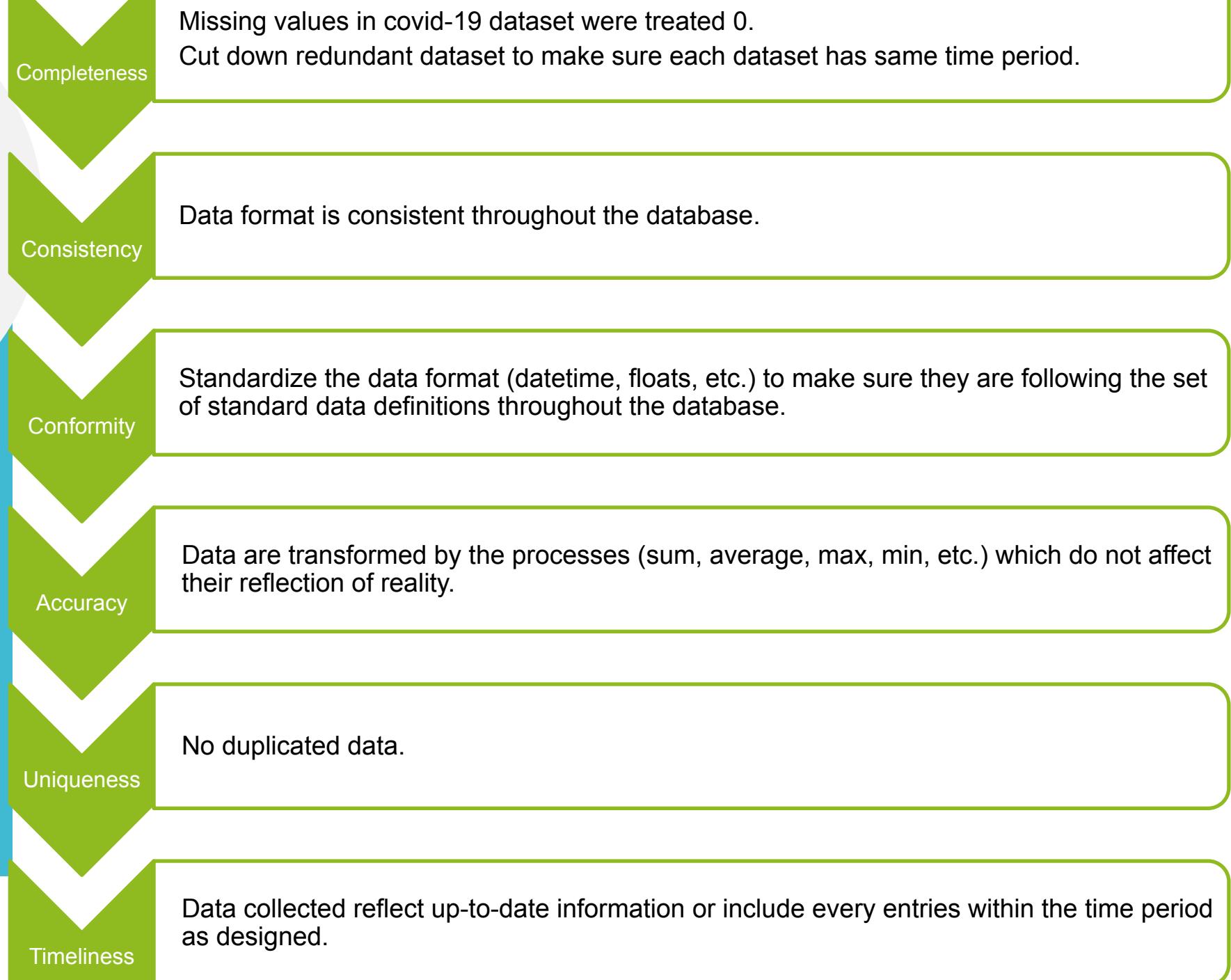
Recommendations



Design Consideration

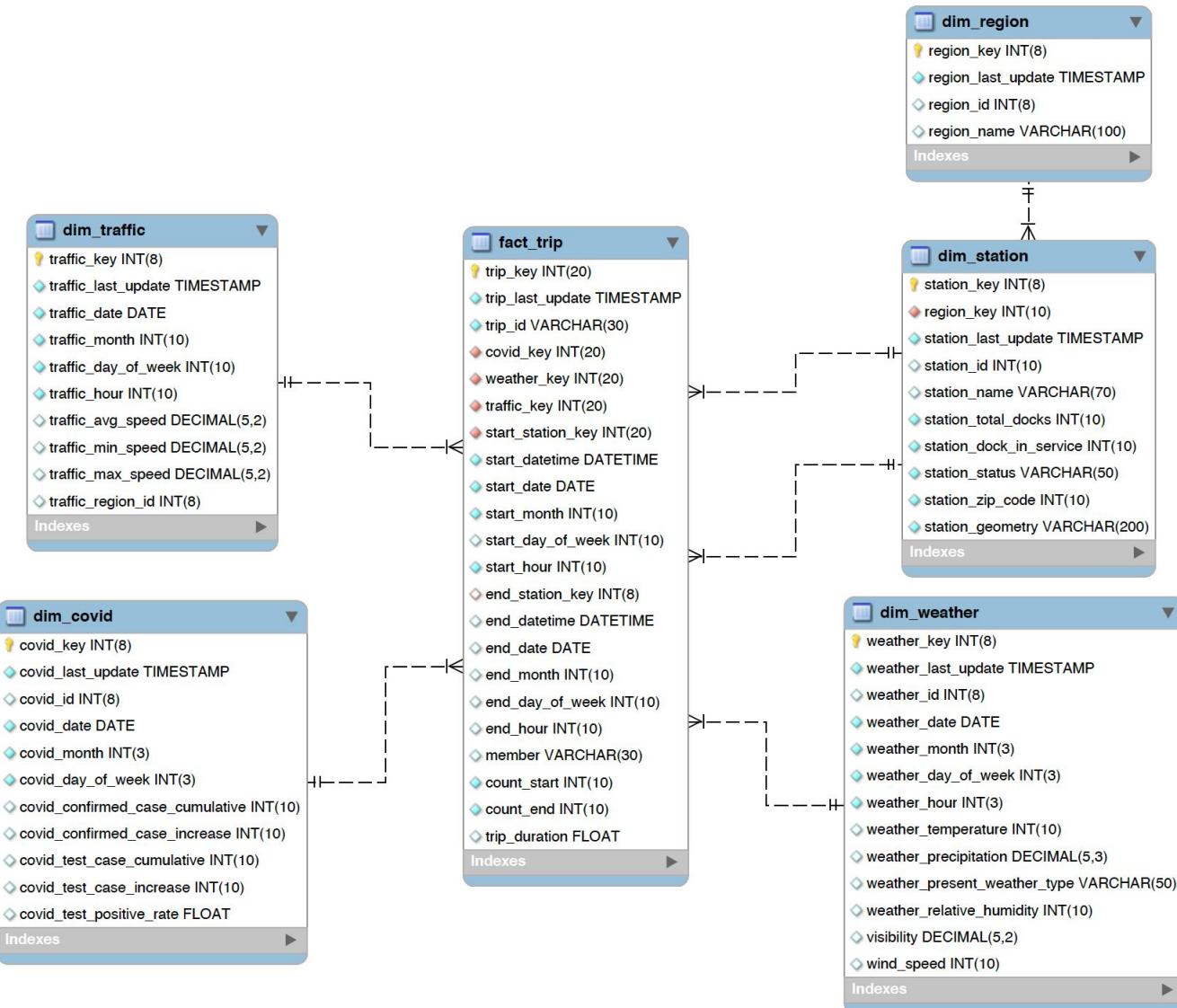
Business Use Case	Data Processing	Data Modeling	Visualization	Recommendations
<h2>Data Types</h2> <ul style="list-style-type: none">• INT: id, key, numbers without decimals• DATE: date• VARCHAR: name, type• Decimal: Other non-integral number	<h2>Dealing with N/As</h2> <ul style="list-style-type: none">• Weather: fill N/A with nearest previous value• Covid: treat N/A with 0• Other dataset: drop row with N/A	<h2>Using Dimensional Table</h2> <ul style="list-style-type: none">• Aggregate and maintain historical data• Optimized for analytical queries• More denormalized and reduces time on processing	<h2>Expected Output of Data Analysis</h2> <ul style="list-style-type: none">• Consumer behavior revealed by data pattern that assists operational decision-making• Relationship between ridership and environmental factors: weather, covid-19, traffic	

Data Quality



Dimensional Model

Linked by Start Time
Or
Linked by Start Region



NoSQL

[MongoDB]



- Easy to manage large amounts of data
- High speed quarries
- Easily scalable
- Limited analytical capacity

Business Use Case Data Processing Data Modeling Visualization Recommendations

localhost ▾ divvybike ▾ 🔒 🛡️ ⚙️ ⏷

```
5
6 → db.TripData.aggregate([
7   {"$sort": {"start_station_id": -1}},
8   {"$group": {"_id": {"start_station_id": "$start_station_id"}, "Count": {"$sum": 1}}}
9 ]))
```

TripData 13.563 s Fetch Count 100 ↕ ⏵ ⏵

Key	Value
▷ { start_station_id : 219 }	{ Count : 1520 }
▷ { start_station_id : 627 }	{ Count : 2340 }
▷ { start_station_id : 158 }	{ Count : 2052 }
▷ { start_station_id : 548 }	{ Count : 26 }
▷ { start_station_id : 408 }	{ Count : 166 }
▷ { start_station_id : 347 }	{ Count : 880 }
▷ { start_station_id : 469 }	{ Count : 154 }
▷ { start_station_id : 591 }	{ Count : 86 }
▷ { start_station_id : 713 }	{ Count : 36 }
▷ { start_station_id : 286 }	{ Count : 928 }

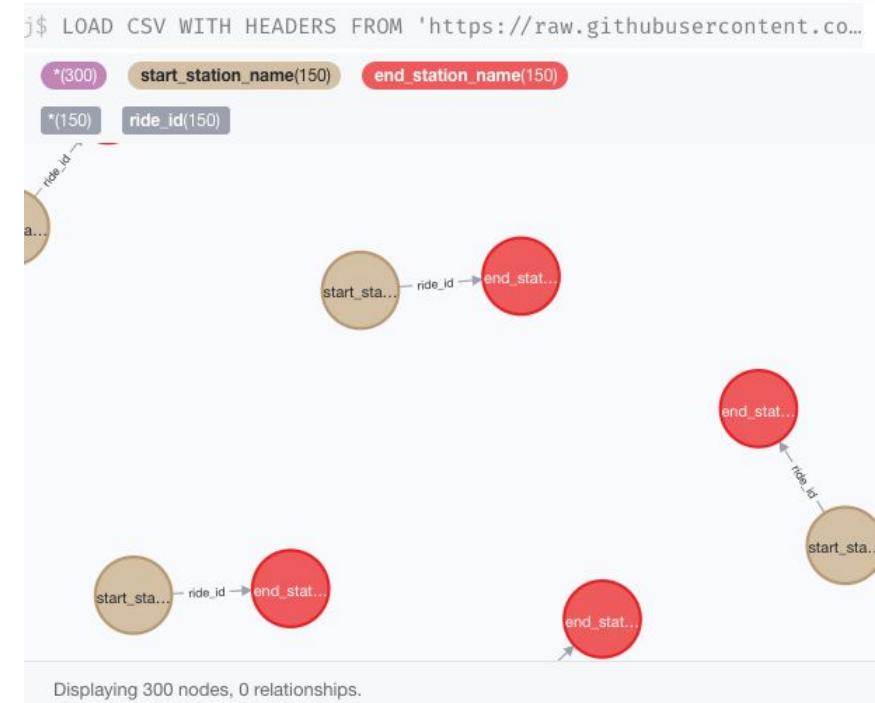
localhost
 ▷ admin
 ▷ config
 ▷ divvybike (6 | 0.13GB)
 ▷ CovidCases (94.0K)
 ▷ DivvyBicycleStations (669)
 ▷ TripData (0.78M | 0.11GB)
 ▷ WeatherData (1.0K)
 ▷ ZipCodes (61)
 ▷ users (0)

NoSQL [Neo4j]



```
1 LOAD CSV WITH HEADERS FROM
  'https://raw.githubusercontent.com/McKenzieKay/DivvyBike/main/202004-divvy-
  tripdata.csv' AS line
2 MERGE (R:ride_id{name:'ride_id'})
3 CREATE (m:start_station_name { name: 'start_station_name' })
4 CREATE (n:end_station_name {name: 'end_station_name' })
5 CREATE (m)-[:ride_id]→(n)
6 RETURN m,n
7
```

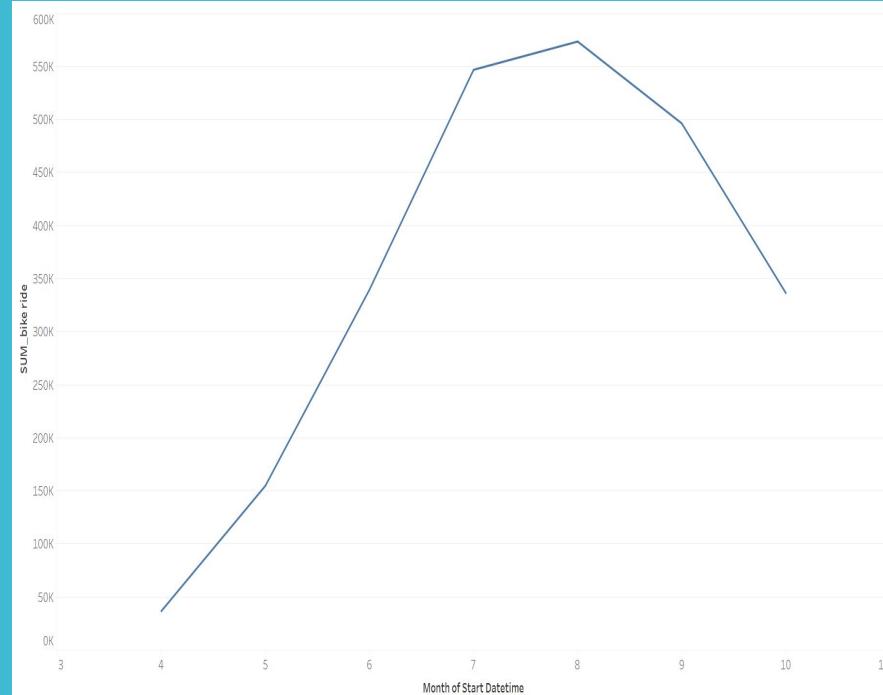
- Had More Troubles
- Created a URL on GitHub to Upload .CSV File
- Still Needs Work on Matching and Creating Relationships



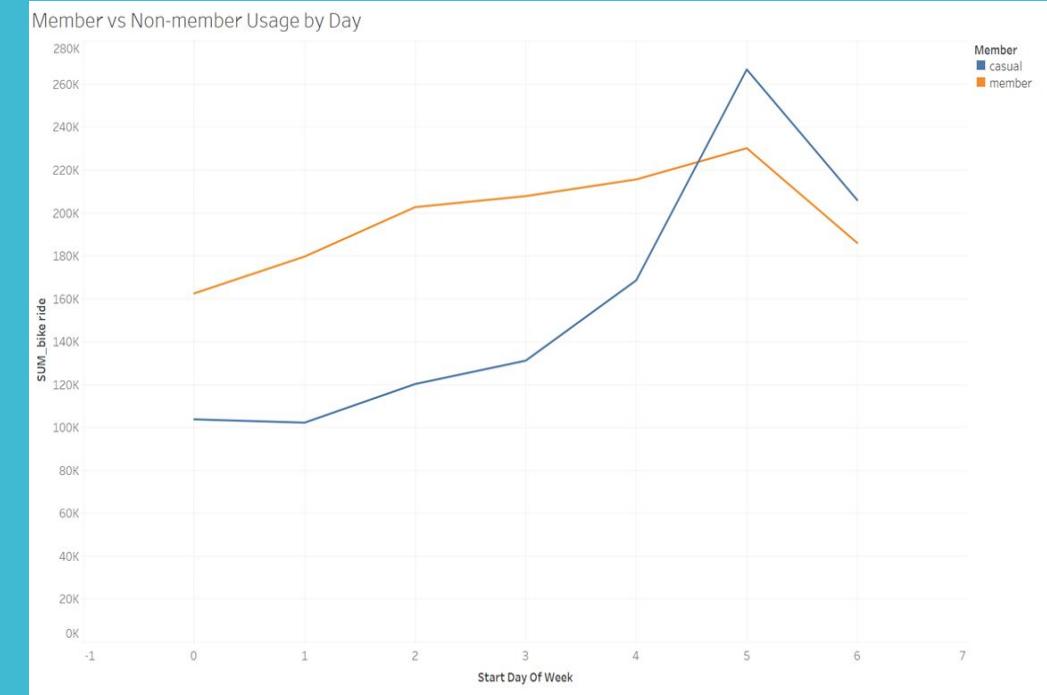
Identifying Factors

- Day of week / Time of Year
- Weather Type/ Temperature
- Covid Cases
- Automobile Traffic

Bike Usage by Month



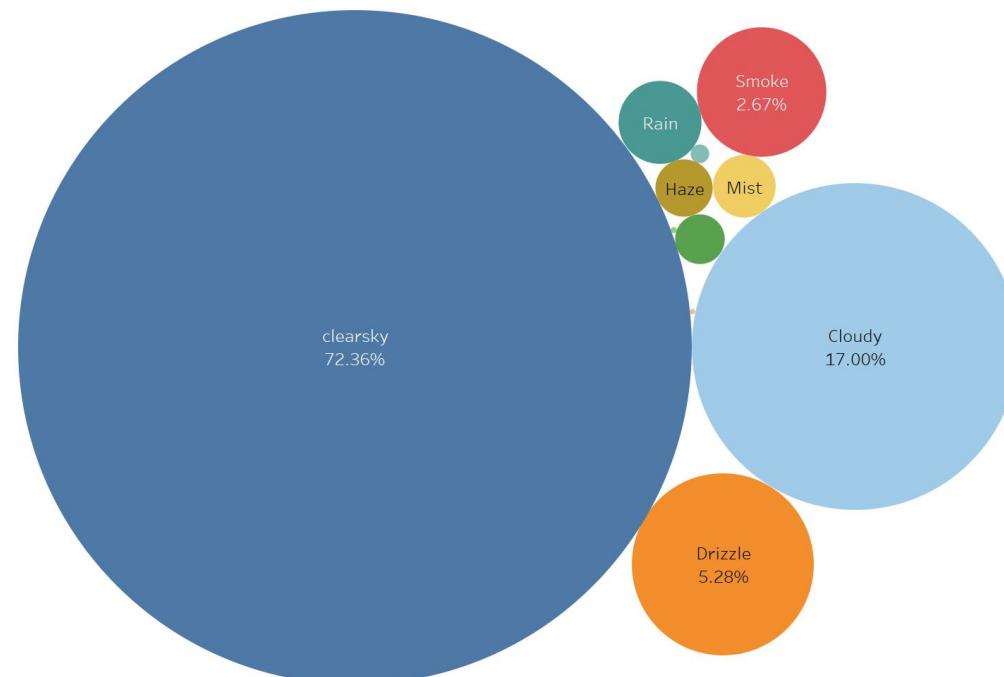
Member vs Non-Member Bike Usage by Day-of-Week



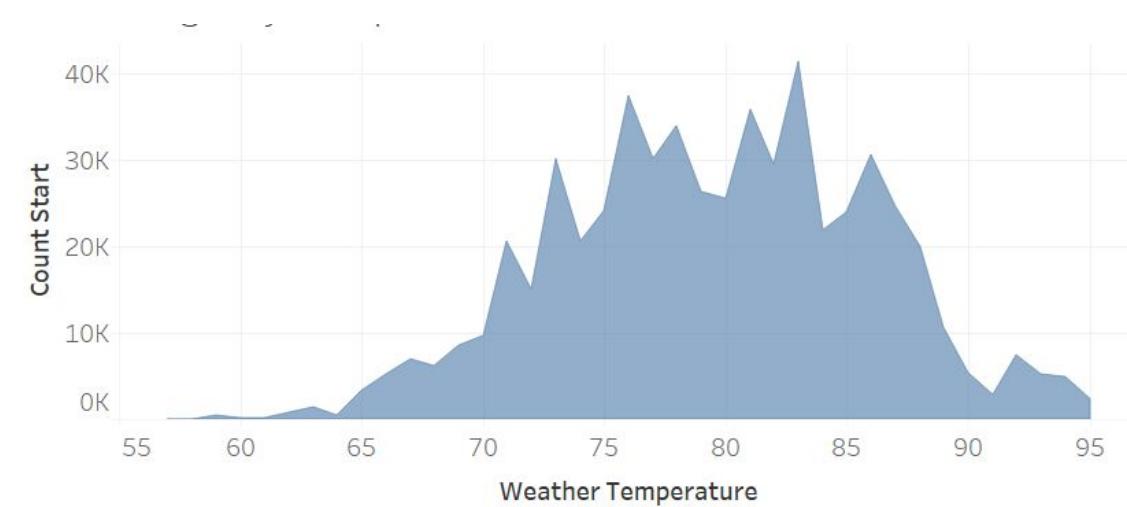
Weather vs Avg Daily Ridership per Hour

Weather..	Apr	May	Jun	Jul	Aug	Sep	Oct	Avg Daily Ridership per..
Clearsky	196.2	285.8	577.7	816.0	809.6	800.7	522.9	
Cloudy	81.6	281.9	347.0	611.6	807.6	657.9	365.7	28.7 940.4
Drizzle	73.5	125.3	404.3	566.5	940.4	429.6	317.3	
Rain	34.1	28.7	340.3	414.9	376.4	190.4	310.0	

Bike Usage by Weather Type

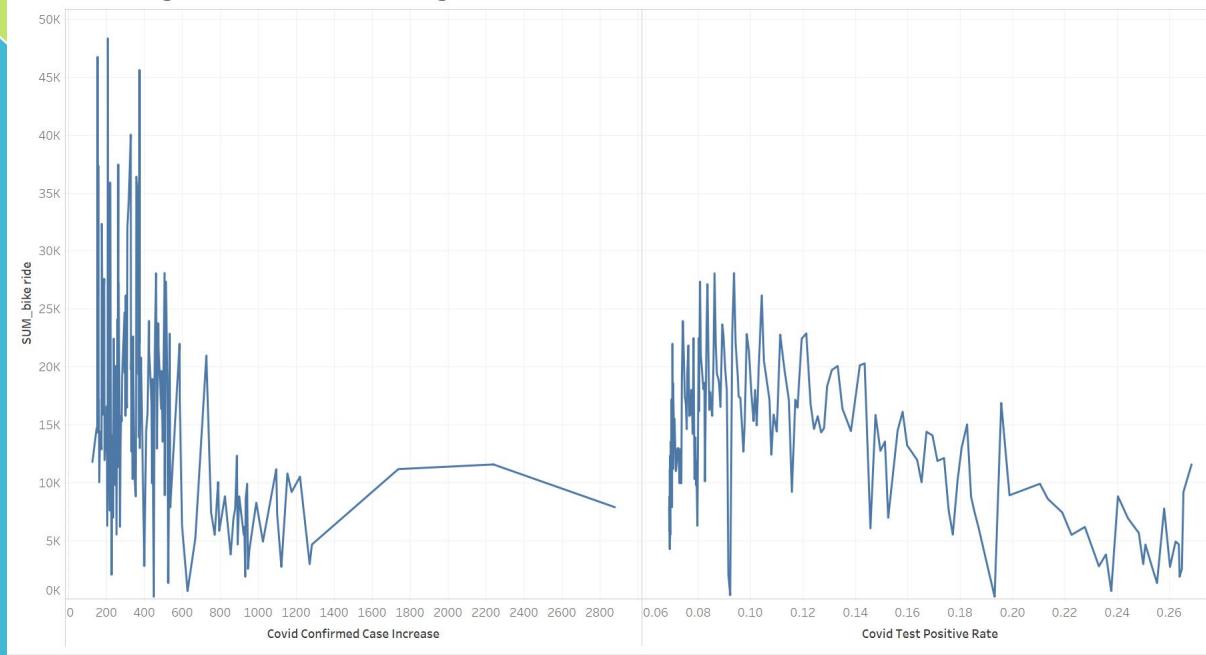


Bike Usage by Temperature



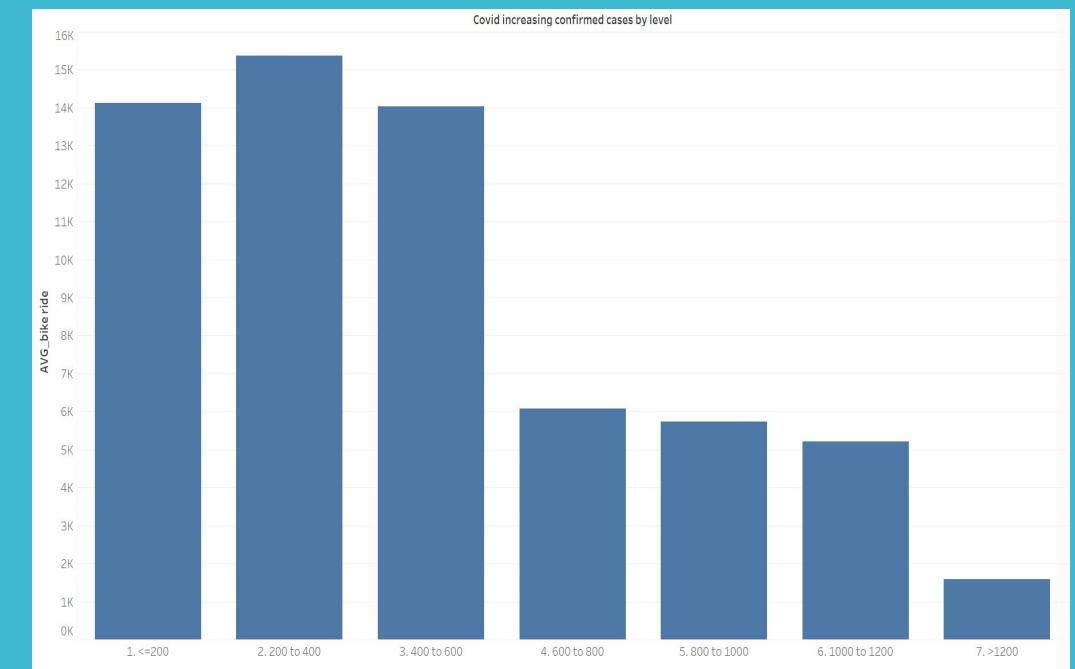
Analysis

Count_Bike rides

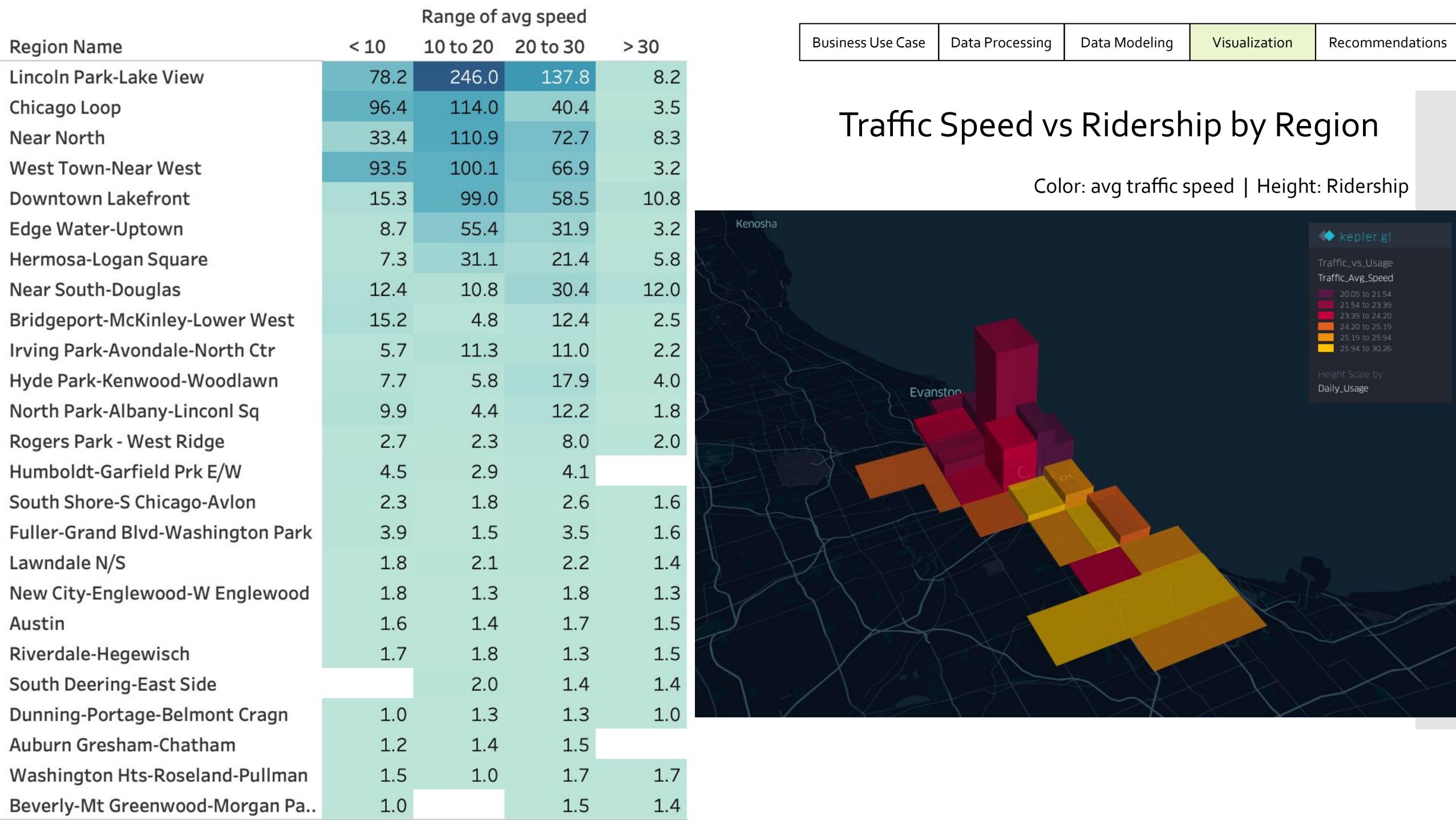


Covid Daily Cases

Covid Daily Test Positive Rate



Covid Daily Increasing Cases by Group vs Avg Bike Rides



Optimizing Allocation

Morning

Traffic is flowing to
'L' stops and
Downtown

Evening

Flowing from
Downtown Out

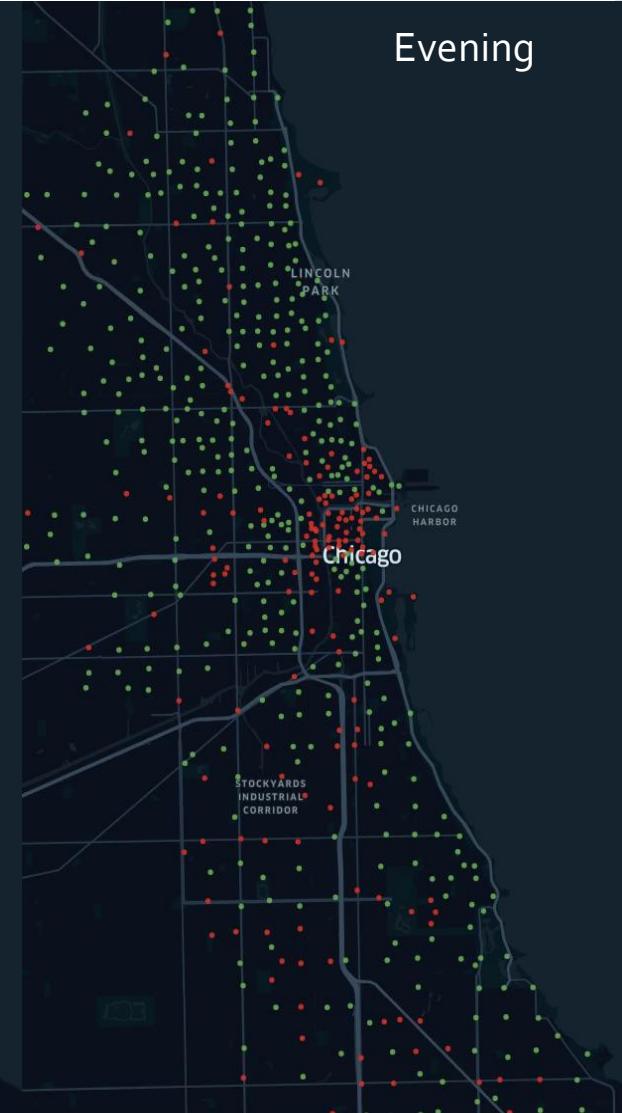
Red: Bike Deficit
Green: Bike Surplus



Morning



Afternoon



Evening

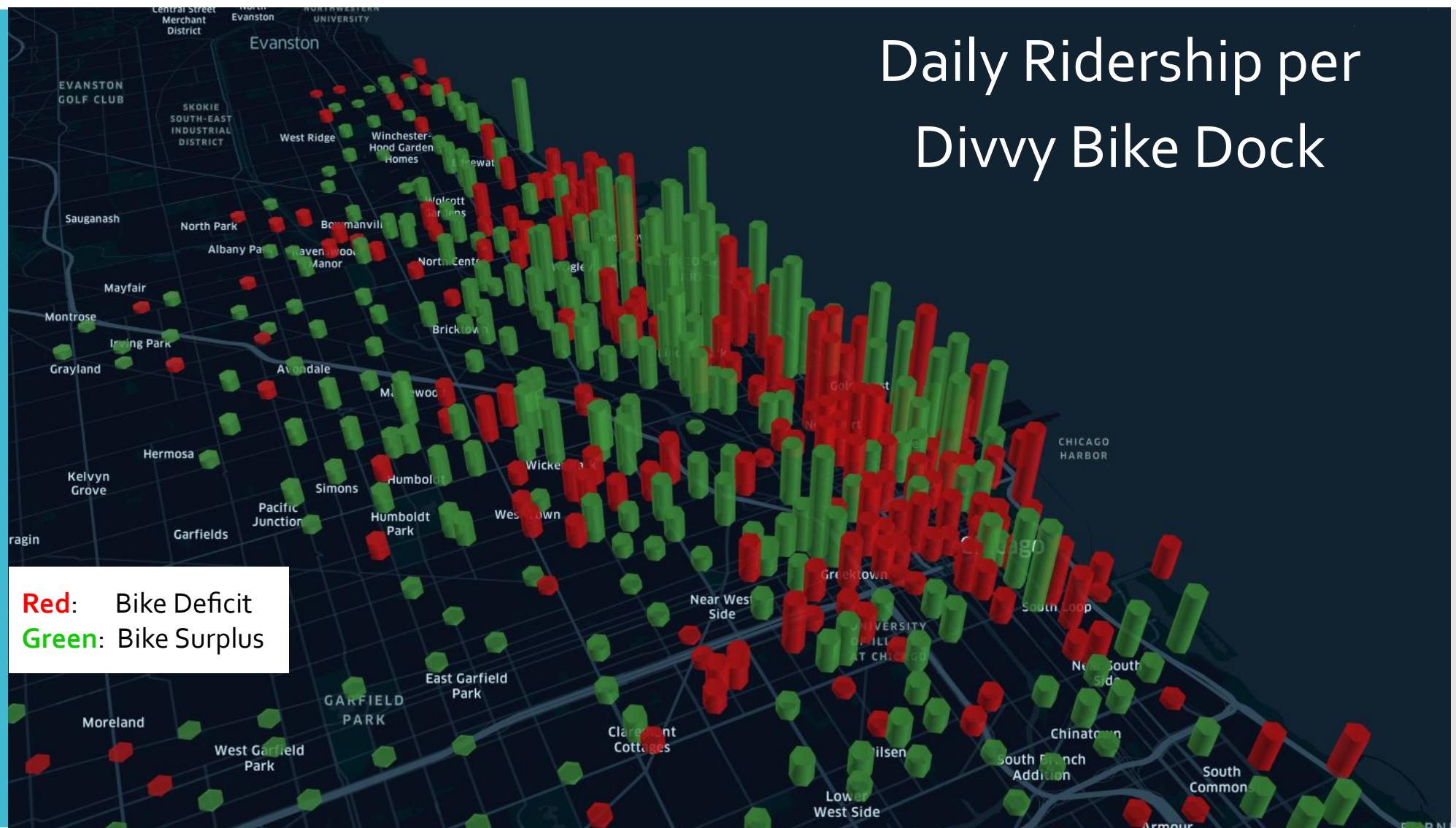
Business Use Case	Data Processing	Data Modeling	Visualization	Recommendations
-------------------	-----------------	---------------	---------------	-----------------

Trip Duration & Ridership by Region



Building New Stations

Add in the vicinity of
TALL **and** RED
column





Recommendations

Identifying Local Factors

Strong Relationship:

- Day of week / Time of Year
- Weather Type / Temperature
- Covid Cases

Weak Relationship:

- Traffic

Optimizing the allocation of bikes

- Morning: More Bikes in Suburban Areas
- Evening: More Bikes in Downtown and 'L' Station

Building new Divvy Stations

Areas with:

- Bike Deficit
- High Ridership

Thank You For Joining Us!

-McKenzie Campbell, Louise Wang, Howard Lin

Appendix

Future Work

- Covid-19 and Weather data with better geographic granularity. (eg. by neighborhood region or by zip code)
- Introduce more informations (demographic data, surrounding public transportation, etc.)
- Better time granularity of covid-19 data will help us find more correlation with ridership
- Visualizing route pattern will provide more information of rider habit

Lesson Learned

- **Problem:** We had problem directly loading CSV files on MySQL.
(LOAD DATA LOCAL INFILE "<filelocation>" INTO TABLE <table name>)
Solution: Import wizard or GCP cloud is used instead.
- **Problem:** Long waiting time importing large data on Mysql shell.
Solution: Upload data into cloud storage and import them directly into cloud sql.
- **Problem:** Original traffic data contains traffic speed in every 10 mins which is far more than we need as by hour.
Solution: *pandas.DataFrame.groupby & pandas.DataFrame.reset_index*

Lesson Learned

- **Problem:** Data are truncated during insert.
Solution: Make sure data format in target field match with selected source.
- **Problem:** Fail to populate large fact table due to small script mistake.
Solution: Try to run the query in smaller dataset first.
- **Problem:** DML with multiple join queries prolong the computation.
Solution: Adding index helps system navigate data and facilitates the matching process.

References

Data Sources



- Divvy Bike historical trip data:
divvzbikes.com
<https://www.divvzbikes.com/system-data>
- Divvy Bike station location:
City of chicago data portal
<https://data.cityofchicago.org/Transportation/Divvy-Bicycle-Stations/bbyy-e7gq>
- Geographic mapping (zip codes boundaries):
City of chicago data portal
<https://data.cityofchicago.org/Facilities-Geographic-Boundaries/Boundaries-ZIP-Codes/gdcf-axmw>
- Weather (Chicago midway weather station):
NOAA data access
<https://www.ncdc.noaa.gov/cdo-web/datatools/lcd>
- Covid-19 weekly cases by zip code:
Illinois Department of Health
[Historical Illinois COVID-19 ZIP code data / COVID-19 in Illinois / Observable \(observablehq.com\)](Historical Illinois COVID-19 ZIP code data / COVID-19 in Illinois / Observable (observablehq.com))
- Chicago Traffic Tracker
City of chicago data portal
<https://data.cityofchicago.org/Transportation/Chicago-Traffic-Tracker-Historical-Congestion-Esti/kf7e-cur8>