



UChicago | MSCA 31012
Data Engineering Platforms for Analytics

**Annie Qurat ul ain
WooJong Choi
Tam Nguyen
Markus Wehr**

Outline

- Executive Summary
- Business Use Case
- Relational database and tools
- Data Analysis and Visualization
- Tableau Visualization
- Summary



Executive Summary



INTRODUCTION

Goal

Invest US\$50 million to:

- Expand stations to all **50 city wards**
- Add **175 stations** and **10,500 bikes**

2019 - 2020



2021

- 2019: More than **20k rides** per day in peak seasons.
- March 2019, **Lyft** took over Divvy
- Early 2020: Plan to pass **20 millionth rides** mark.

Second expansion
(107 new stations)

Provided its 15
millionth rides in 2018

2015 - 2016



2017 - 2018

First expansion
(175 new stations)

Officially launched
in June 2013
(75 stations and 750
bikes)



2013



Bikeshare system



6,000 bikes



608 stations

***Chicagoans' regular mode
of transportation***

RESEARCH OBJECTIVES

- To assist with the expansion plan, our team developed a relational database that will enable quick response and analysis on the current state Divvy operations in regard to ridership, station locations and various other factors affecting them. And:

• Provide methodologies and various tools used in the process

• Provide data analysis and visualization

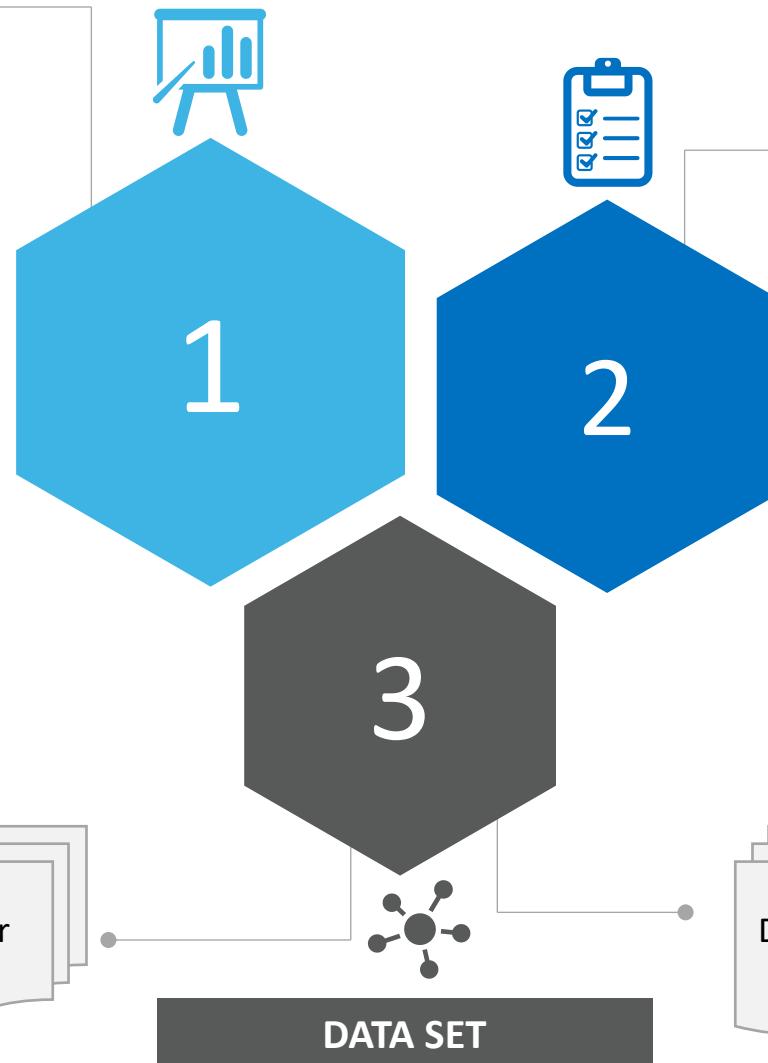
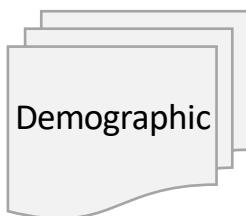
• Put forward a future state blueprint for the new stations and bikes allocation process



PROPOSED FINDING

Our final deliverables will enable Divvy leadership to:

- Understand current ridership and station locations
- Understand various factors that impact ridership. i.e
 - Demographic
 - Traffic volume
 - Bike racks / lanes
 - Weather
- Develop dashboards and KPIs to gauge overall business / operation performance
- Plan for future station & bikes allocation



METHODOLOGY

- Develop a scoring model to determine optimal number of stations and bikes by zip codes based on various factors
- Visualize findings from analysis - trends, outliers, patterns and predictions



Data Source



Dataset	Source		File Format	Size
Trip	Divvy	https://www.divvybikes.com/system-data	CSV	> 1mil rows
Station	City of Chicago	https://data.cityofchicago.org/Transportation/Divvy-Bicycle-Stations/bbyy-e7gq	CSV	> 600 rows
Station_zip	Divvy	https://feeds.divvybikes.com/stations/stations.json	JSON	> 600 rows
Weather	National Weather Service Forecast Office	https://w2.weather.gov/climate/xmacis.php?wfo=lot	CSV	> 12k rows
Bike racks	City of Chicago	https://data.cityofchicago.org/Transportation/Bike-Racks/cbyb-69xx	CSV	> 5k rows
Population	City of Chicago	https://catalog.data.gov/dataset?res_format=CSV&organization=city-of-chicago	CSV	< 100 rows
Bike route	City of Chicago	https://data.cityofchicago.org/Transportation/Bike-Routes/3w5d-sru8	CSV	< 1k rows
Zip code	Chicago Data Type	http://robparal.blogspot.com/2013/07/chicago-community-area-and-zip-code.html	CSV	< 100 rows



Relational Database and Tools

Fact and dimensional table



Table Name	Table Type	Cardinality	Additional Details
fact_trip	Fact Table	M:1 Relationship with Station and Weather Table	Contains information about each trip including the start/end station, total time, age, gender of the customer
dim_station	Dimensional Table	1:M relationship with Fact Table	Contains information like station address, total number of docks available, date the station became available.
dim_weather	Dimensional Table	1:M relationship with Fact Table	Contains temperature, rain/snow, wind information in hourly format. Also, contains the sunset and sunrise time.
dim_population	Dimensional Table	1:M relationship with Location Table	Contains information about the population (age, gender) demographics zip wise.
dim_location	Dimensional Table	M:1 relationship with Population Table	Contains the location of all the stations, traffic routes, bike routes. Zip code is a must have for each address.
dim_traffic	Dimensional Table	1:M relationship with Location Table	Contains the traffic flow information daily including the direction (Northbound, Southbound, Westward, Eastward) on streets.
dim_bike_racks	Dimensional Table	1:M relationship with Location Table	Contains information about the non-divvy bike racks scattered across Chicago city
dim_bike_lane	Dimensional Table	1:1 relationship with Location Table	Contains information about the bike routes in the city, including their length and the streets they run on.

Fact table joined with Dimension tables provides interesting insights into how variables interact. Fact Table can be sliced by time and diced by stations, gender and age variables.

Database Design: Enhanced Entity Relational Diagram



Dimensional Schema: SNOWFLAKE

DDL

```

-- PGOOL_Workbench Forward Engineering

SET GLOBAL validate_checksums=0;
SET GLOBAL unique_checks=0;
SET GLOBAL FOREIGN_KEY_CHECKS=0;
SET GLOBAL SQL_MODE='NO_ZERO_DATE,NO_ZERO_IN_DATE,NO_ZERO_DATE_ERROR,FOREIGN_KEY_CHECKS=0';

-- Schema divey
CREATE SCHEMA IF NOT EXISTS divey DEFAULT CHARACTER SET utf8;
USE divey;
-- Table: divey.dim_weather
CREATE TABLE IF NOT EXISTS `divey.dim_weather` (
    `station_id` INT NOT NULL,
    `name` VARCHAR(255) NOT NULL,
    `temperature` INT NOT NULL,
    `wind` INT NOT NULL,
    `precipitation` INT NOT NULL,
    PRIMARY KEY (`name`)
) ENGINE=InnoDB;

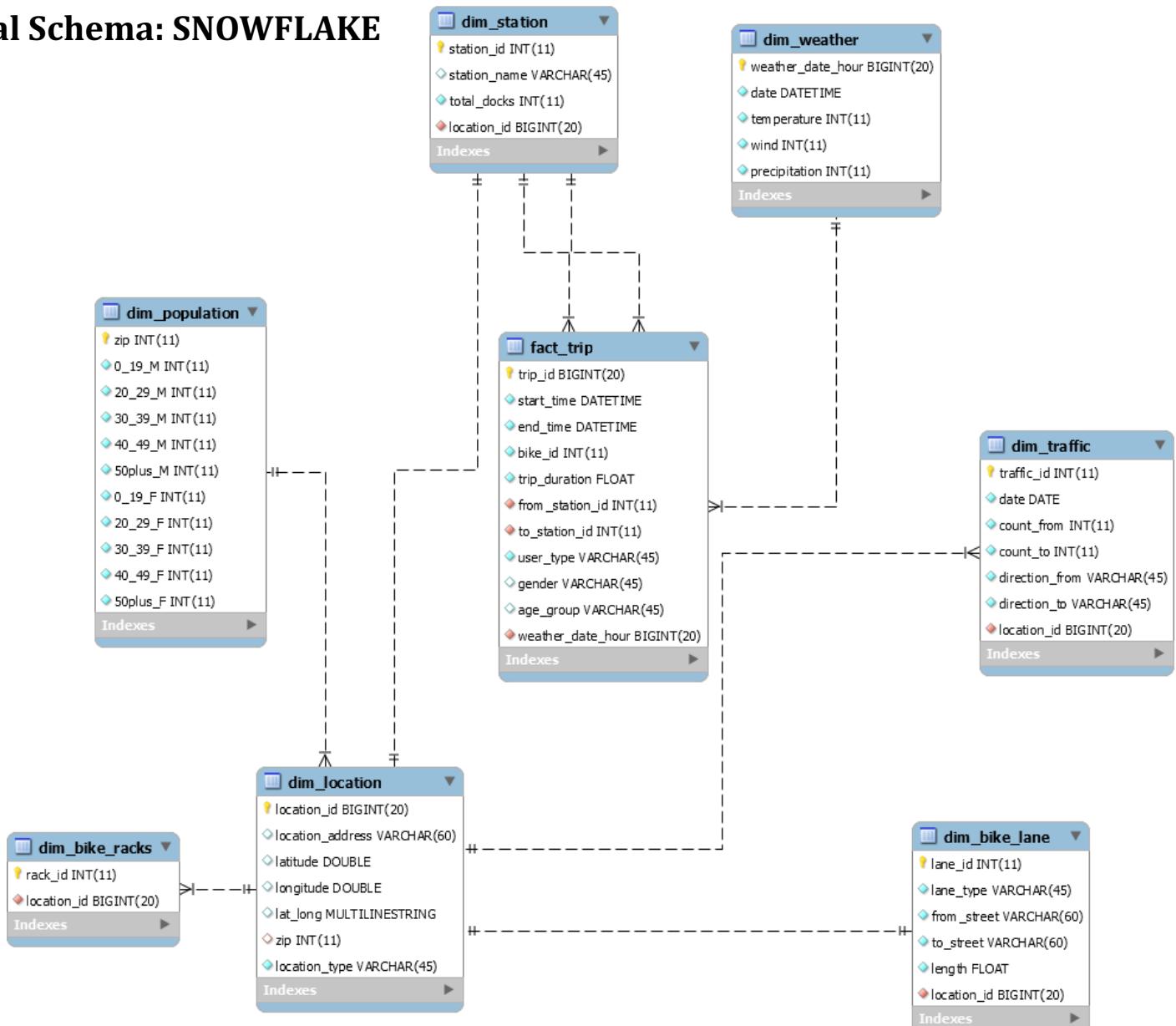
-- Table: divey.dim_location
CREATE TABLE IF NOT EXISTS `divey.dim_location` (
    `location_id` INT NOT NULL,
    `name` VARCHAR(255) NOT NULL,
    `lat` DECIMAL(10, 8) NOT NULL,
    `lon` DECIMAL(10, 8) NOT NULL,
    `elevation` INT NOT NULL,
    PRIMARY KEY (`name`)
) ENGINE=InnoDB;

-- Create table if not exists divey.dim_station
CREATE TABLE IF NOT EXISTS `divey.dim_station` (
    `station_id` INT NOT NULL,
    `name` VARCHAR(255) NOT NULL,
    `lat` DECIMAL(10, 8) NOT NULL,
    `lon` DECIMAL(10, 8) NOT NULL,
    `elevation` INT NOT NULL,
    `location_id` INT NOT NULL,
    `CONSTRAINT `location_id` FOREIGN KEY(`location_id`) REFERENCES `divey.dim_location`(`location_id`),
    PRIMARY KEY (`name`),
    FOREIGN KEY(`location_id`) REFERENCES `divey.dim_location`(`location_id`)
) ENGINE=InnoDB;

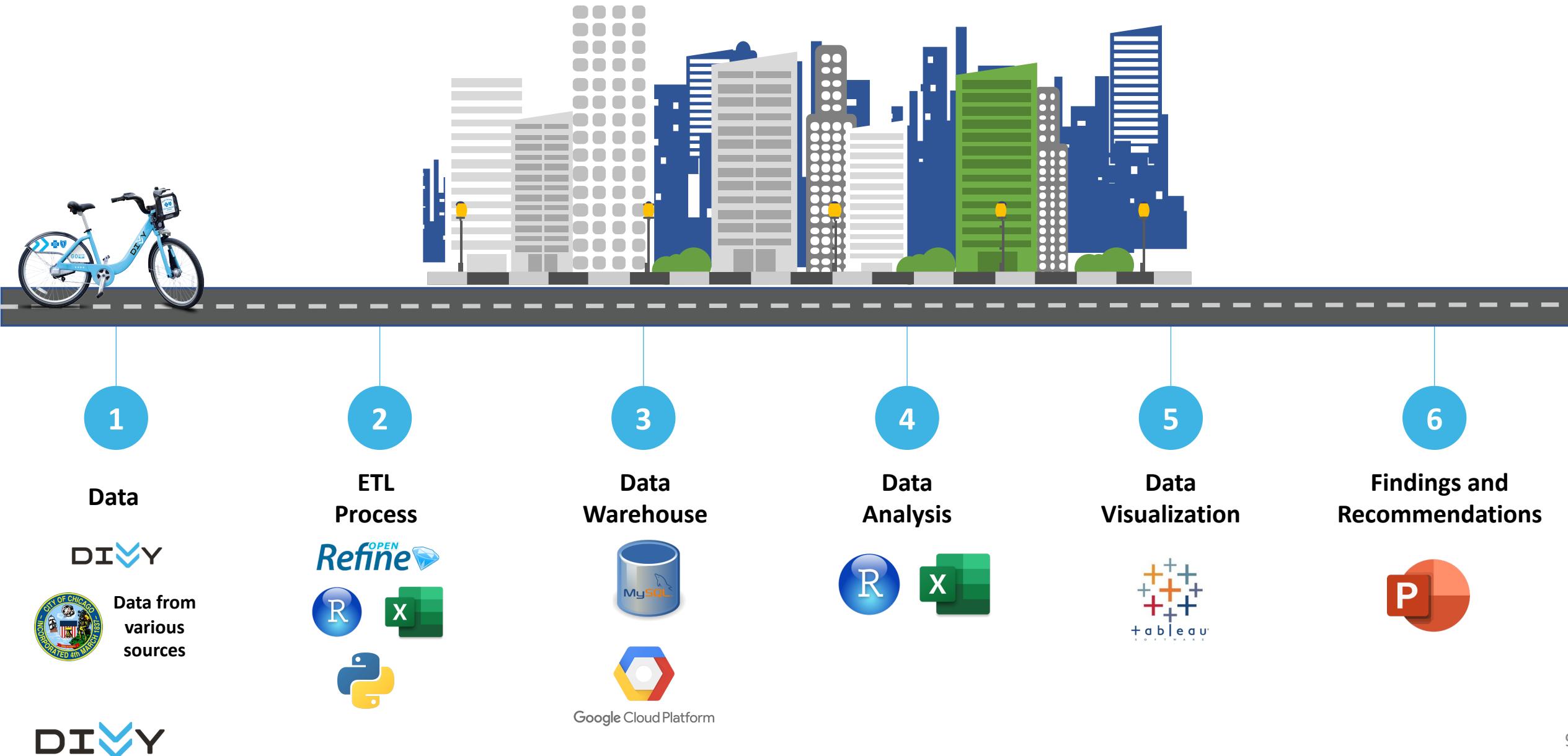
-- Table: divey.dim_traffic
CREATE TABLE IF NOT EXISTS `divey.dim_traffic` (
    `traffic_id` INT NOT NULL,
    `count_from` INT NOT NULL,
    `count_to` INT NOT NULL,
    `count_over` INT NOT NULL,
    `percentage` DECIMAL(10, 2) NOT NULL,
    `CONSTRAINT `location_id` FOREIGN KEY(`location_id`) REFERENCES `divey.dim_location`(`location_id`)
) ENGINE=InnoDB;

```

DML



Tools



Data extraction, Cleaning, Normalization



- Create and load database
 - Produce queries to support project's analysis purpose



- Import, clean, and extract real-time station data from Divvy to get the zip code for each station.



- Clean all dimensional tables to import to mySQL
 - Analyze descriptive data: customer profiling, zip, stations
 - Build the scoring system for research objectives' purpose: add more stations and bikes.

- Get the zipcode using longitude and latitude for dim_location table
 - Estimated the distance between trips
 - Stack the distance data to produce an adaptable format for tableau visualization purpose
 - Conduct some correlation between trips and other factors: weekday, bike racks, weather...

```
BIKE RACKS, weather...
```

```
1 install.packages("rvegeo")
2 library(rvegeo)
3
4
5 dataPath <- "/Users/tommy/Tommy/UChicago/Data_Engineering_Platform/FinalProject/Processing/"
6
7 #rock1
8 rock1 <- read.csv(paste0(dataPath,"rock1.csv",sep = ""))
9 rock1_zip <- rvegeo::longitude_rock1(zipcode, latitude=rock1$latitude,
10                                     longitude=rock1$longitude, type="zip")
11 write.table(rock1_zip, file = paste0(dataPath,"rock1.zip.csv",sep = ''), row.names = F)
12
13 #rock2
14 rock2 <- read.csv(paste0(dataPath,"rock2.csv",sep = ""))
15 rock2_zip <- rvegeo::longitude_rock2(zipcode, latitude=rock2$latitude,
16                                     longitude=rock2$longitude, type="zip")
17 write.table(rock2_zip, file = paste0(dataPath,"rock2.zip.csv",sep = ''), row.names = F)
18
19 #rock3
20 rock3 <- read.csv(paste0(dataPath,"rock3.csv",sep = ""))
21 rock3_zip <- rvegeo::longitude_rock3(zipcode, latitude=rock3$latitude,
22                                     longitude=rock3$longitude, type="zip")
23 write.table(rock3_zip, file = paste0(dataPath,"rock3.zip.csv",sep = ''), row.names = F)
```

```
10 # distance
11 library(general)
12 #distance<-CO
13 #for(i in 1:50)
14 #  distance<-c(distance,distVIncentyEllipsoid(c(lat1,lon1),c(lat2,lon2)))
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
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
```



- Construct fact_trip table to import to my SQL:
 - Calculate the age group of Divvy users
 - Add in new column as a foreign key using in my

```
In [34]: def age_group(age):
    if age < 20:
        return '0-19'
    else:
        age -= 20
        if age < 10:
            return '20-29'
        else:
            age -= 10
            if age < 10:
                return '30-39'
            else:
                age -= 10
                if age < 10:
                    return '40-49'
                else:
                    age -= 10
                    if age < 10:
                        return '50-59'
                    else:
                        age -= 10
                        if age < 10:
                            return '60-69'
                        else:
                            age -= 10
                            if age < 10:
                                return '70-79'
                            else:
                                age -= 10
                                if age < 10:
                                    return '80+'

fact_trip[fact_trip['age_group'] = fact_trip['age'].apply(age_group)

fact_trip[fact_trip['age_group'].value_counts()>=10]

Out[34]: 20-29    245732
          30-39    13842
          40-49    12777
          50-59    126001
          60-69      177
          70-79      101
          80+       101
          Name: age_group, dtype: int64

In [35]: print(fact_trip.head(10))

trip_id          start_time        end_time     bike_id trip_duration \
1  1000000000 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
2  1000000001 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
3  1000000002 2018-04-01 00:11:44 2018-04-01 00:12:19   5145    945.0
4  1000000003 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
5  1000000004 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
6  1000000005 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
7  1000000006 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
8  1000000007 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
9  1000000008 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
10 1000000009 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0

from_station_id to_station_id  user_type  gender  birthyear  age
0           100         100  Subscriber   Male  1953.0  54.0
1           157         100  Subscriber   Male  1953.0  54.0
2           149         100  Subscriber   Male  1953.0  54.0
3           241         171  Subscriber   Male  1959.0  21.0
4           149         100  Subscriber   Male  1953.0  54.0
5           149         100  Subscriber   Male  1953.0  54.0
6           149         100  Subscriber   Male  1953.0  54.0
7           149         100  Subscriber   Male  1953.0  54.0
8           149         100  Subscriber   Male  1953.0  54.0
9           149         100  Subscriber   Male  1953.0  54.0

In [37]: fact_trip['date'] = pd.DatetimeIndex(fact_trip['start_time']).date
fact_trip['month'] = pd.DatetimeIndex(fact_trip['start_time']).month
fact_trip['day'] = pd.DatetimeIndex(fact_trip['start_time']).day
fact_trip['hour'] = pd.DatetimeIndex(fact_trip['start_time']).hour
fact_trip['weekofyear'] = pd.DatetimeIndex(fact_trip['start_time']).isocalendar().week
fact_trip['weather_desc_hour'] = fact_trip['weather_desc'].map(lambda x: x[x.find(')')+1:])

/Users/amy/anaconda/lib/python3.6/site-packages/pandas/core/settings.py:310: SettingWithCopyWarning:
A value is being coerced to integer type and may lose precision.
Try using .loc[[row_index], col_index] = value instead

See the caveats in the documentation:
http://pandas.pydata.org/pandas-docs/stable/indexing.html#coercing-view-references
ValueError: invalid literal for int() with base 10: 'valley'

Out[37]: fact_trip[fact_trip['date'].dt.year == 2018].head(10)

trip_id          start_time        end_time     bike_id trip_duration \
0  1000000000 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
1  1000000001 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
2  1000000002 2018-04-01 00:11:44 2018-04-01 00:12:19   5145    945.0
3  1000000003 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
4  1000000004 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
5  1000000005 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
6  1000000006 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
7  1000000007 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
8  1000000008 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
9  1000000009 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0

fact_trip[fact_trip['date'].dt.month == 4].head(10)

trip_id          start_time        end_time     bike_id trip_duration \
0  1000000000 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
1  1000000001 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
2  1000000002 2018-04-01 00:11:44 2018-04-01 00:12:19   5145    945.0
3  1000000003 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
4  1000000004 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
5  1000000005 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
6  1000000006 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
7  1000000007 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
8  1000000008 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
9  1000000009 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0

fact_trip[fact_trip['date'].dt.day == 1].head(10)

trip_id          start_time        end_time     bike_id trip_duration \
0  1000000000 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
1  1000000001 2018-04-01 00:11:44 2018-04-01 00:12:03   3619    499.0
2  1000000002 2018-04-01 00:11:44 2018-04-01 00:12:19   5145    945.0
3  1000000003 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
4  1000000004 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
5  1000000005 2018-04-01 00:11:44 2018-04-01 00:12:23   5065    945.0
6  1000000006 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
7  1000000007 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
8  1000000008 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
9  1000000009 2018-04-01 00:11:44 2018-04-01 00:12:23   3619    499.0
```

Sample Queries



Net influx per station and hour

```
5 •   SELECT
6     TripFrom.station_id,
7     TripFrom.stationName AS stationName,
8     TripFrom.TimeOfDay AS tripTime,
9     TripFrom.tripFrom,
10    TripTo.tripTo,
11    (TripFrom.tripFrom - TripTo.tripTo) AS NetTrip
12
13   FROM
14
15   (SELECT
16     ds.station_id,
17     ds.station_name AS stationName,
18     ds.total_docks AS totalDocks,
19     HOUR(ft.start_time) AS TimeOfDay,
20     COUNT(ft.from_station_id) as tripFrom
21
22   FROM
23     fact_trip ft
24       INNER JOIN
25       dim_station ds ON ds.station_id = ft.from_station_id
26
27   GROUP BY
28     ds.station_id, TimeOfDay
29
30   ORDER BY
31     ds.station_id, TimeOfDay ASC) AS TripFrom
32
33   INNER JOIN
34
35   (SELECT
36     ds.station_id,
37     ds.station_name AS stationName,
38     ds.total_docks AS totalDocks,
39     HOUR(ft.end_time) AS TimeOfDay,
40     COUNT(ft.to_station_id) as tripTo
41
42   FROM
43     fact_trip ft
44       INNER JOIN
45       dim_station ds ON ds.station_id = ft.to_station_id
46
47   GROUP BY
48     ds.station_id, TimeOfDay
49
50   ORDER BY ds.station_id, TimeOfDay ASC) AS TripTo ON TripFrom.station_id = TripTo.station_id
51
52 WHERE TripFrom.TimeOfDay = TripTo.TimeOfDay;
```

Average distance travelled per station and zip code

```
91 •   SELECT
92
93     FrS.station_id,
94     FrS.trip_id,
95     FrS.latitude AS lat1,
96     FrS.longitude AS long1,
97     TrS.station_id,
98     TrS.trip_id,
99     TrS.latitude AS lat2,
100    TrS.longitude AS long2
101
102   FROM
103
104   (SELECT
105     ds.station_id,
106     ft.trip_id,
107     dl.latitude,
108     dl.longitude
109
110   FROM
111     dim_location dl
112       INNER JOIN
113       dim_station ds ON dl.location_id=ds.location_id
114
115   INNER JOIN
116     fact_trip ft ON ds.station_id=ft.from_station_id) AS FrS
117
118   INNER JOIN
119
120   (SELECT
121     ds.station_id,
122     ft.trip_id,
123     dl.latitude,
124     dl.longitude
125
126   FROM
127     dim_location dl
128       INNER JOIN
129       dim_station ds ON dl.location_id=ds.location_id
130
131   INNER JOIN
132     fact_trip ft ON ds.station_id=ft.to_station_id) AS TrS ON FrS.trip_id=TrS.trip_id
133
134 WHERE
135     FrS.station_id != TrS.station_id;
```

Data Analysis and Visualization



Customer Profiling

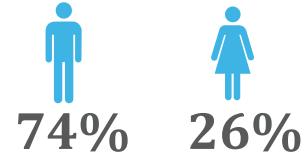


Users Type

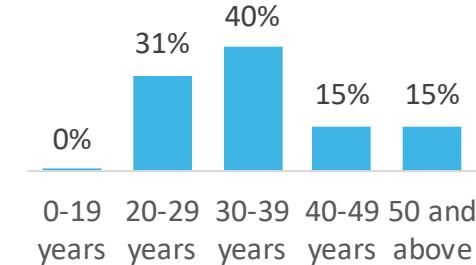


■ Subscriber
■ Non-subscriber

Gender



Age Group



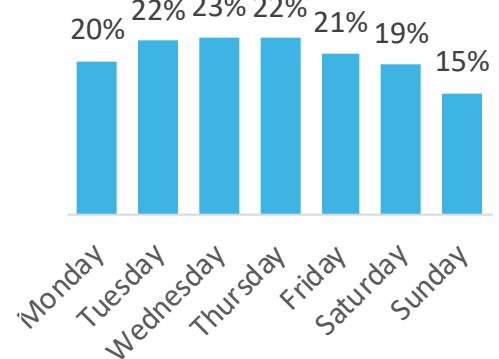
Average Distance Travel



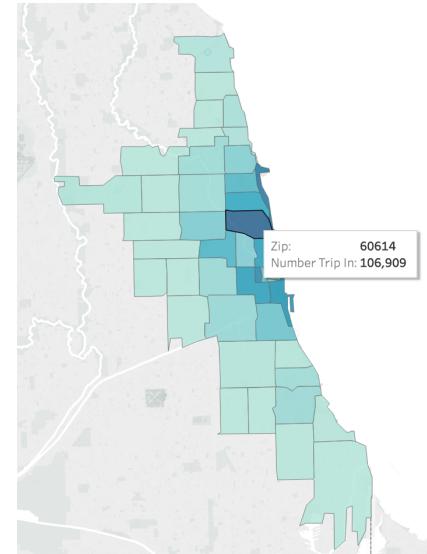
0.95 miles



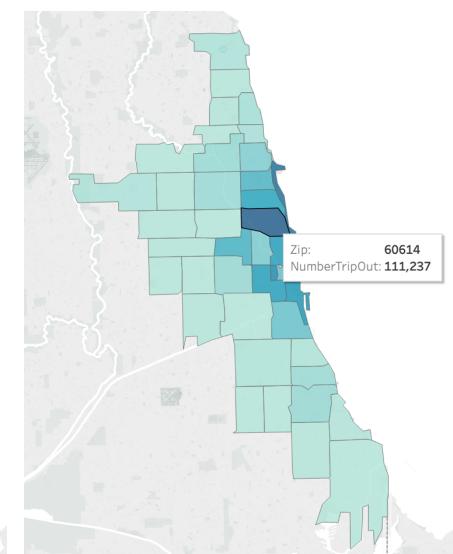
Trip by day



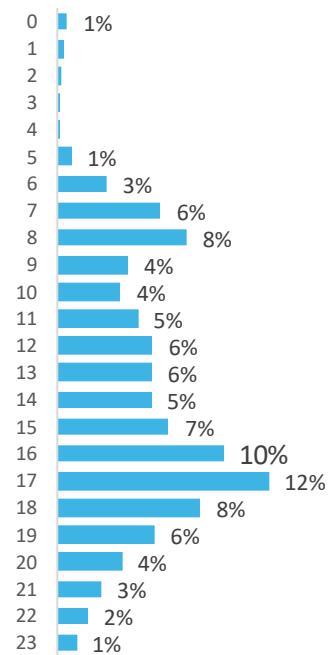
Trip-in by area



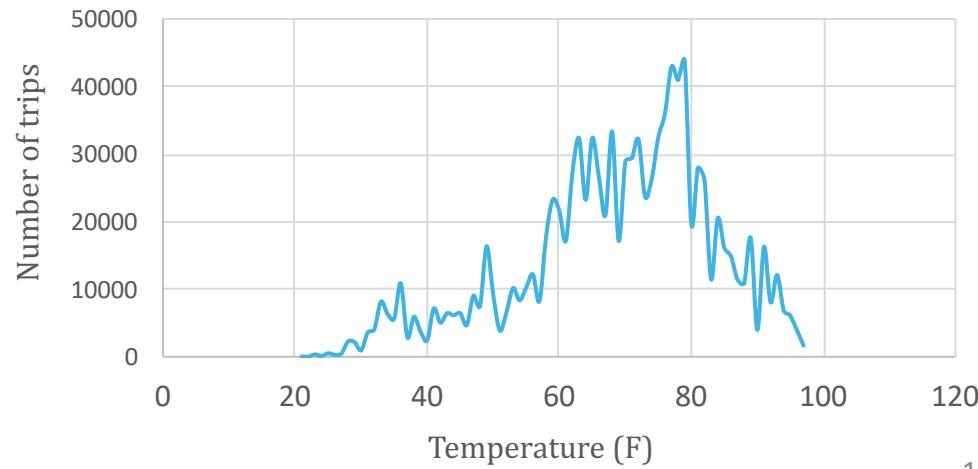
Trip-out by area



Trip by hour



Trip by weather



Findings by zip code



Zipcode Analysis

Population	Gender	Demographic		Traffic	DCK	Number of Bike racks	Divvy Stations		Divvy Trips																
		Age					Number of stations	Total # of docks	Avg # of docks	Avg of Avg Distance from other stations (miles)		Trips Out	Trips In	Net											
		Total	Male%	Female%	0_19	20_29	30_39	40_49	50plus	0_19	20_29	30_39	40_49	50plus											
60605	23,594	48,545	51.5%	0.0678	0.138	0.0394	0.0585	0.0134	0.0719	0.1451	0.163	0.0603	0.0206	8,300	356	28.9	68,302	65,243	(3,659)	6.33%	60.7%				
60601	1115	48,454	50.6%	0.0659	0.0702	0.093	0.0585	0.0103	0.0716	0.1744	0.0953	0.0571	0.0206	23,800	1	33.1	4,59	68,506	63,594	(4,912)	6.23%	72.4%			
60609	6436	50,454	50.4%	0.0659	0.0702	0.093	0.0585	0.0103	0.0716	0.1744	0.0953	0.0571	0.0206	18,875	75	28.9	34,302	32,462	(1,840)	6.04%	61.9%				
60649	13,654	48,524	51.17%	0.0671	0.0653	0.0549	0.0584	0.0152	0.0789	0.073	0.0689	0.054	0.0206	19,300	31	14	170	12.1	3,79	1,522	1,482	(40)	6.04%	62.0%	
60614	66623	47,674	52.4%	0.0571	0.0642	0.0686	0.0527	0.0951	0.0627	0.181	0.0375	0.051	0.043	21,700	195	34	638	18.8	4,82	106,935	111,237	4,328	10.23%	73.3%	
60608	82743	47,674	47.6%	0.1505	0.0604	0.0569	0.0587	0.1355	0.0368	0.0968	0.0616	0.0503	0.0204	18,600	109	27	370	13.7	5.03	14,489	15,288	759	1.4%	67.8%	
60622	52,593	51.5%	48.9%	0.0682	0.1414	0.0354	0.0593	0.0666	0.0645	0.1352	0.1295	0.0566	0.0631	34,400	354	25	465	18.6	4,57	45,14	46,260	146	4.3%	66.0%	
60609	23,594	48,545	51.5%	0.0678	0.138	0.0394	0.0585	0.0134	0.0719	0.1451	0.163	0.0603	0.0206	8,300	356	28.9	68,302	65,243	(3,659)	6.33%	60.7%				
60607	23,594	48,545	51.5%	0.0671	0.0687	0.0545	0.0587	0.0137	0.0622	0.0906	0.0893	0.0536	0.0176	29,300	62	26	453	17.4	4,40	61,385	61,617	232	5.8%	31.1%	
60642	18,665	51.5%	48.9%	0.0688	0.1424	0.0361	0.0593	0.0130	0.0656	0.084	0.1100	0.9	0.021	21,700	9	12	217	18.1	4,37	24,446	24,523	107	2.3%	90.8%	
60610	37730	48,454	52.4%	0.0397	0.1339	0.0395	0.0533	0.0456	0.0545	0.0547	0.1057	0.0616	0.0362	22,100	122	18	403	14.5	4,52	61,694	62,462	708	5.8%	81.3%	
60624	28722	48,454	52.4%	0.0394	0.1334	0.0375	0.0533	0.0456	0.0458	0.0545	0.1059	0.0616	0.0362	22,100	149	16	233	23.3	4,67	83,207	95,133	9,326	7.0%	58.4%	
60604	47,674	48,545	50.8%	0.0708	0.0589	0.0504	0.0586	0.0713	0.0738	0.0587	0.0522	0.0113	0.0206	11,700	11	4	558	10.5	4,62	63,705	63,231	(534)	88.3%	60.7%	
60603	497	48,324	50.7%	0.0704	0.163	0.0526	0.0589	0.0107	0.0724	0.1751	0.0946	0.0522	0.0127	13,700	98	5	135	20.7	4,56	27,542	24,721	(2,821)	2.47%	60.7%	
60616	47,674	48,545	51.5%	0.0658	0.0598	0.0494	0.0576	0.0622	0.0493	0.0368	0.0684	0.0526	0.0204	6,100	87	29	447	15.4	5.35	32,058	32,470	323	3.05%	60.7%	
60602	7,730	48,524	50.7%	0.0622	0.0646	0.0553	0.0587	0.0624	0.0606	0.0693	0.0652	0.0543	0.0204	13,700	32	12	354	15.5	4,34	78,901	76,461	(2,440)	16.2%	60.7%	
60637	49,006	48,324	51.5%	0.0682	0.0646	0.0553	0.0587	0.0624	0.0606	0.0693	0.0652	0.0543	0.0204	10,900	72	17	264	12.5	8,15	12,584	12,429	(25)	1.1%	81.1%	
60657	66601	48,75	50.3%	0.0525	0.0662	0.0153	0.0585	0.0532	0.0562	0.068	0.0565	0.0685	0.0204	12,900	201	20	371	16.6	5.28	54,456	57,328	2,872	5.26%	60.9%	
60647	82,397	50,454	50.7%	0.0511	0.0612	0.0172	0.053	0.039	0.0608	0.0591	0.062	0.0606	0.0204	10,600	200	15	333	14.5	4,03	26,516	27,351	835	2.47%	60.7%	
60612	47,674	48,545	51.5%	0.0658	0.0598	0.0553	0.0587	0.0624	0.0606	0.0693	0.0652	0.0543	0.0204	13,700	100	15	333	14.5	4,03	26,516	27,351	835	2.47%	60.7%	
60615	40,008	48,324	51.5%	0.0591	0.0674	0.0563	0.0587	0.0624	0.0606	0.0693	0.0652	0.0543	0.0204	10,600	64	11	177	16.1	7,54	10,823	10,307	94	1.03%	76.7%	
60618	32,085	50,25	48.9%	0.1238	0.0955	0.0107	0.0703	0.0595	0.0637	0.0597	0.1007	0.1223	0.0596	0.0204	18,700	193	21	308	14.7	5.95	10,373	11,778	806	1.07%	85.2%
60613	48,295	50,454	48.6%	0.0502	0.164	0.0164	0.0568	0.0626	0.0108	0.0593	0.0625	0.0139	0.0204	11,600	76	23	426	18.5	5.90	43,431	44,295	864	4.14%	79.5%	
60625	23,592	48,524	50.5%	0.0504	0.1613	0.0162	0.0563	0.0624	0.0107	0.0592	0.0625	0.0139	0.0204	13,700	83	13	262	13.3	6.62	41,101	42,274	176	1.3%	73.3%	
60640	65,736	51,852	48.2%	0.0775	0.1103	0.0783	0.0399	0.0727	0.0303	0.0332	0.0228	0.0104	0.0204	13,700	205	14	282	20.1	6.71	22,613	23,060	447	2.16%	80.4%	
60619	63,630	48,351	48.5%	0.0577	0.0595	0.0572	0.0585	0.0642	0.074	0.0674	0.0642	0.0704	0.0204	17,800	31	17	183	10.8	9,86	6,58	684	46	0.06%	77.6%	
60634	74,302	43.7	50.3%	0.1233	0.0697	0.0711	0.0589	0.089	0.1273	0.0721	0.0736	0.0708	0.0204	21,800	112	1	19	15.0	5.33	3,164	3,653	523	0.32%	73.8%	
60620	47,674	50,454	49.6%	0.0651	0.0654	0.0505	0.0543	0.0677	0.0664	0.0623	0.0671	0.0651	0.0204	14,700	145	15	325	8.2	8,12	7,25	7,255	60	0.44%	60.7%	
60621	35,918	48,324	50.8%	0.0517	0.0654	0.0505	0.0543	0.0677	0.0664	0.0623	0.0671	0.0651	0.0204	15,000	28	13	130	15.8	8,28	10,321	10,321	(14)	0.04%	88.3%	
60624	38,703	48,624	53.4%	0.051	0.063	0.0607	0.0589	0.0653	0.0644	0.0744	0.0634	0.0604	0.0204	11,800	35	8	88	11.0	5.71	3,330	320	(19)	0.03%	68.1%	
60623	32,121	53.7%	46.4%	0.1844	0.0991	0.0601	0.0634	0.0732	0.0673	0.0653	0.0737	0.0661	0.0543	0.0204	13,700	109	6	66	11.0	5.88	560	588	28	0.05%	82.0%
60645	45,288	48,524	50.5%	0.1342	0.0763	0.0686	0.0608	0.0404	0.0369	0.0774	0.0696	0.0533	0.0204	11,400	55	7	161	12.1	8,99	1,236	1,237	1	7.22%	60.7%	
60630	47,674	48,545	51.5%	0.0671	0.0686	0.0608	0.0686	0.0502	0.0753	0.0634	0.0744	0.0652	0.0204	14,700	50	4	44	7.9	7,93	10,507	10,507	0	0.05%	65.5%	
60660	47,521	48.2%	51.8%	0.0844	0.0403	0.0563	0.0791	0.0606	0.0802	0.0512	0.0753	0.0634	0.0204	35,200	42	5	31	18.2	7.57	5,215	5,686	481	0.1%	79.7%	
60641	7,736	48,524	50.4%	0.0709	0.0803	0.0687	0.0612	0.1377	0.0802	0.0875	0.0708	0.1358	0.0204	39,900	81	6	78	13.0	7,21	1,173	1,110	(63)	0.1%	77.2%	
60630	54,093	48,251	50.3%	0.1221	0.0795	0.0768	0.0757	0.0499	0.1259	0.0733	0.0789	0.0741	0.0653	0.0204	12,300	37	2	26	13.0	7,28	230	239	9	0.02%	73.5%
60651	49,552	48,752	50.7%	0.0674	0.0724	0.0598	0.0585	0.0624	0.0624	0.0661	0.0674	0.0684	0.0204	14,700	11	6	65	12.5	6,25	1,504	1,795	243	0.02%	60.7%	
60644	49,552	48,752	50.7%	0.0674	0.0724	0.0598	0.0585	0.0624	0.0624	0.0661	0.0674	0.0684	0.0204	14,700	11	6	65	12.5	6,25	1,504	1,795	243	0.02%	60.7%	
60636	49,552	48.624	53.4%	0.1587	0.0671	0.0527	0.0587	0.0818	0.0785	0.0603	0.0873	0.0673	0.0204	18,200	20	8	80	10.0	8,10	105	109	4	0.01%	76.2%	
60637	49,552	48.624	53.4%	0.1587	0.0671	0.0527	0.0587	0.0818	0.0785	0.0603	0.0873	0.0673	0.0204	18											

Score based approach



Current station locations (Before expansion plan)

Where are the stations?

- CTA, Metra stations
- employment centers, shopping districts, medical centers, schools
- other popular destinations.

How were the locations chosen?

- population density
- business permits
- other stations in the surrounding network.

Our scoring methodology

- When Divvy first launched, it focused more on the popular destinations (tourist attraction areas, shopping centers, offices etc.)
- The expansion plan is focused more on expanding to the areas where there are currently no Divvy stations
- Priority = underserved communities (in terms of number of Divvy stations).
- Score based system for the allocation of the stations and the bikes taking into consideration the below factors. New station allocation determined based on overall score (i.e. higher score = more stations)

Category	Score Description	Weight	Comments
Divvy Stations (existing)	less number of stations = more points	↓ ↑	20% More weight assigned to zip codes with no stations. Points deducted to zip codes with stations
Trips (Trips Out)	more number of trips = more points	↑ ↑	10% -
Net (Trip From - Trip To)	lower Net value = more points	↓ ↑	5% Points only added to zip codes with a negative net value
Subscriber%	higher % of subscribers = more points	↑ ↑	15% -
Population Total	higher population = more points	↑ ↑	15% -
Male%	higher male % = more points	↑ ↑	5% -
20_39 Age Group	higher % of 20_39 age group = more points	↑ ↑	10% -
Average Distance to other stations	higher avg distance to other stations = more points	↑ ↑	10% -
Traffic	higher vehicle volume = more points	↑ ↑	5% -
Bike racks	more number of bike racks (bike friendliness score) = more points	↑ ↑	5% -

Scores by zip code



Scoring by zip calculation

Demographic	Traffic		Bike Racks		Divvy Stations		Divvy Trips						Population	Gender	20_39		Vehicle Volume	Number of Bike racks	Number of stations	Area of Avg distance from other stations (mi)	Trips (Trips Out - Trip To)	Net (Trip From - Trip To)	Subscriber %																							
	Total	Male%	20k	10k	10k	5k	10k	10k	5k	10k	10k	5k	10k	10k	5k	20k	10k	10k	5k																											
	Total	Male%	20k	10k	10k	5k	10k	10k	5k	10k	10k	5k	10k	10k	5k	20k	10k	10k	5k																											
Weight	15%	5%	10%	5%	5%	20%	10%	10%	5%	10%	10%	5%	10%	10%	5%	20%	10%	10%	5%	Total Points	1500	500	1000	500	500	2000	1000	1000	500	1500	1500	500	1500	1500	500	2000	1000	1000	500							
50605	15072	48.5%	0.934	8,300	356	19	4.75	68,902	(3,659)	60.7%	12.7	8.7	23.9	3.9	(3,25)	15.4	65.0	84.3	24.4	240.0	5	308	15072	48.5%	0.934	8,300	356	19	4.75	68,902	(3,659)	60.7%	12.7	8.7	23.9	3.9	(3,25)	15.4	65.0	84.3	24.4	240.0	5	308		
50601	1115	49.4%	0.523	23,800	1	11	4.59	88,506	(4,912)	72.4%	6.2	8.8	25.0	11.1	0.1	(18.09)	14.3	64.7	103.2	23.2	255.0	5	326	1115	49.4%	0.523	23,800	1	11	4.59	88,506	(4,912)	72.4%	6.2	8.8	25.0	11.1	0.1	(18.09)	14.3	64.7	103.2	23.2	255.0	5	326
50609	14591	50.2%	0.366	8,100	75	28	6.49	29,953	36.1	9.0	14.4	3.8	6.7	46.05	21.1	24	-	35.2	62.6	2	106	14591	50.2%	0.366	8,100	75	28	6.49	29,953	36.1	9.0	14.4	3.8	6.7	46.05	21.1	24	-	35.2	62.6	2	106				
50649	48654	43.8%	0.263	18,300	31	H	9.79	1,922	(60)	62.6%	26.0	7.8	12.6	8.6	2.8	(23.03)	31.8	14	14	25.0	94.3	2	121	48654	43.8%	0.263	18,300	31	H	9.79	1,922	(60)	62.6%	26.0	7.8	12.6	8.6	2.8	(23.03)	31.8	14	14	25.0	94.3	2	121
50616	66623	47.6%	0.310	21,700	195	34	4.82	106,909	4,328	79.2%	37.1	8.5	24.9	10.2	17.5	(55.92)	15.7	100.3	-	31.9	190.7	4	245	66623	47.6%	0.310	21,700	195	34	4.82	106,909	4,328	79.2%	37.1	8.5	24.9	10.2	17.5	(55.92)	15.7	100.3	-	31.9	190.7	4	245
50608	82743	52.4%	0.379	18,800	109	27	5.03	14,489	799	87.8%	46.0	9.4	17.8	8.8	9.8	(44.41)	16.3	13.7	-	35.4	112.7	2	145	82743	52.4%	0.379	18,800	109	27	5.03	14,489	799	87.8%	46.0	9.4	17.8	8.8	9.8	(44.41)	16.3	13.7	-	35.4	112.7	2	145
50622	62565	51.5%	0.546	34,400	354	25	4.57	45,114	146	86.0%	29.2	9.1	25.4	16.1	31.7	(41.12)	14.8	42.6	-	34.7	162.6	3	209	62565	51.5%	0.546	34,400	354	25	4.57	45,114	146	86.0%	29.2	9.1	25.4	16.1	31.7	(41.12)	14.8	42.6	-	34.7	162.6	3	209
50606	2314	48.2%	0.528	10,200	227	6	4.40	37,508	(3,420)	90.9%	13	8.8	25.0	4.8	20.3	(3.87)	14.3	35.4	78.8	36.6	215.4	5	277	2314	48.2%	0.528	10,200	227	6	4.40	37,508	(3,420)	90.9%	13	8.8	25.0	4.8	20.3	(3.87)	14.3	35.4	78.8	36.6	215.4	5	277
50607	23902	49.2%	0.533	29,300	62	26	4.40	61,385	232	91.6%	13.3	8.8	25.0	13.7	5.6	(42.76)	14.3	57.9	-	36.7	132.5	3	170	23902	49.2%	0.533	29,300	62	26	4.40	61,385	232	91.6%	13.3	8.8	25.0	13.7	5.6	(42.76)	14.3	57.9	-	36.7	132.5	3	170
50642	18485	51.5%	0.542	11,000	9	12	4.37	24,416	107	90.8%	10.3	9.1	25.6	5.2	0.8	(19.74)	14.2	23.0	-	36.6	105.1	2	135	18485	51.5%	0.542	11,000	9	12	4.37	24,416	107	90.8%	10.3	9.1	25.6	5.2	0.8	(19.74)	14.2	23.0	-	36.6	105.1	2	135
50610	37730	46.4%	0.495	22,100	122	18	4.52	16,594	708	81.9%	210	8.3	22.8	10.3	10.9	(23.61)	14.7	58.2	-	33.0	149.6	3	192	37730	46.4%	0.495	22,100	122	18	4.52	16,594	708	81.9%	210	8.3	22.8	10.3	10.9	(23.61)	14.7	58.2	-	33.0	149.6	3	192
50611	28722	48.4%	0.495	18,100	149	16	4.67	83,207	5,362	58.4%	16.0	8.3	22.8	8.5	13.3	(26.32)	15.2	84.2	-	23.5	165.4	4	213	28722	48.4%	0.495	18,100	149	16	4.67	83,207	5,362	58.4%	16.0	8.3	22.8	8.5	13.3	(26.32)	15.2	84.2	-	23.5	165.4	4	213
50654	14890	48.7%	0.489	23,000	115	14	4.40	48,187	(2,806)	88.3%	8.3	8.8	22.9	10.8	10.3	(23.03)	14.8	64.4	64.6	35.6	216.7	5	278	14890	48.7%	0.489	23,000	115	14	4.40	48,187	(2,806)	88.3%	8.3	8.8	22.9	10.8	10.3	(23.03)	14.8	64.4	64.6	35.6	216.7	5	278
50604	575	49.4%	0.533	11,000	227	3	4.55	15,680	(450)	80.4%	0.3	8.8	25.1	5.1	20.3	(4.93)	14.8	14.8	10.4	32.4	127.1	3	163	575	49.4%	0.533	11,000	227	3	4.55	15,680	(450)	80.4%	0.3	8.8	25.1	5.1	20.3	(4.93)	14.8	14.8	10.4	32.4	127.1	3	163
50603	497	49.3%	0.513	13,700	38	5	4.56	27,542	(2,821)	55.5%	0.3	8.8	24.9	6.4	8.8	(3.22)	14.8	26.0	65.0	22.4	163.2	4	217	497	49.3%	0.513	13,700	38	5	4.56	27,542	(2,821)	55.5%	0.3	8.8	24.9	6.4	8.8	(3.22)	14.8	26.0	65.0	22.4	163.2	4	217
50616	48437	48.1%	0.364	6,100	87	29	5.35	32,058	412	77.6%	27.0	8.6	16.3	7.8	(47.70)	17.4	30.3	-	31.2	94.3	2	121	48437	48.1%	0.364	6,100	87	29	5.35	32,058	412	77.6%	27.0	8.6	16.3	7.8	(47.70)	17.4	30.3	-	31.2	94.3	2	121		
50602	1210	49.3%	0.532	18,750	50	3	4.53	17,152	(2,21)	79.2%	0.7	8.8	25.0	8.8	4.4	(4.93)	14.7	16.2	0.5	31.9	106.1	2	136	1210	49.3%	0.532	18,750	50	3	4.53	17,152	(2,21)	79.2%	0.7	8.8	25.0	8.8	4.4	(4.93)	14.7	16.2	0.5	31.9	106.1	2	136
50651	7298	48.2%	0.531	15,600	102	12	4.24	78,947	(2,434)	93.2%	4.3	8.8	25.0	7.3	8.1	(19.74)	14.1	74.4	56.1	37.6	217.0	5	279	7298	48.2%	0.531	15,600	102	12	4.24	78,947	(2,434)	93.2%	4.3	8.8	25.0	7.3	8.1	(19.74)	14.1	74.4	56.1	37.6	217.0	5	279
50637	48008	44.9%	0.307	20,900	72	17	8.15	12,684	(1,255)	81.4%	27.5	8.0	14.4	6.4	9.8	(3.87)	15.5	28.5	5.9	32.9	115.4	2	148	48008	44.9%	0.307	20,900	72	17	8.15	12,684	(1,255)	81.4%	27.5	8.0	14.4	6.4	9.8	(3.87)	15.5	28.5	5.9	32.9	115.4	2	148
50640	65736	51.8%	0.428	16,700	205	14	6.71	22,613	447	80.4%	36.6	3.3	20.1	7.8	18.4	(23.03)	21.8	21.3	-	32.4	144.6	3	186	65736	51.8%	0.428	16,700	205	14	6.71	22,613	447	80.4%	36.6	3.3	20.1	7.8	18.4	(23.03)	21.8	21.3	-	32.4	144.6	3	186
50619	63830	43.6%	0.234	17,800	31	17	9.86	638	46	77.6%	35.5	7.8	11.0	8.3	2.8	(27.36)	17.6	32.0	0.6	31.2	101.3	2	130	63830	43.6%	0.234	17,800	31	17	9.86	638	46	77.6%	35.5	7.8	11.0	8.3	2.8	(27.36)	17.6	32.0	0.6	31.2	101.3	2	130
50634	74302	48.1%	0.285	21,800	112	1	5.33	1,364	529	79.2%	8.8	8.8	13.4	10.2	10.0	(164)	17.3	3.0	-	32.1	134.6	3	173	74302	48.1%	0.285	21,800	112	1	5.33	1,364	529	79.2%	8.8	8.8	13.4	10.2	10.0	(164)	17.3	3.0	-	32.1	134.6	3	173
50626	50144	50.4%	0.419	7,100	145	15	8.92	7,235	60	94.5%	27.9	9.0	18.9	3.3	13.0	(13.0)	14.7																													

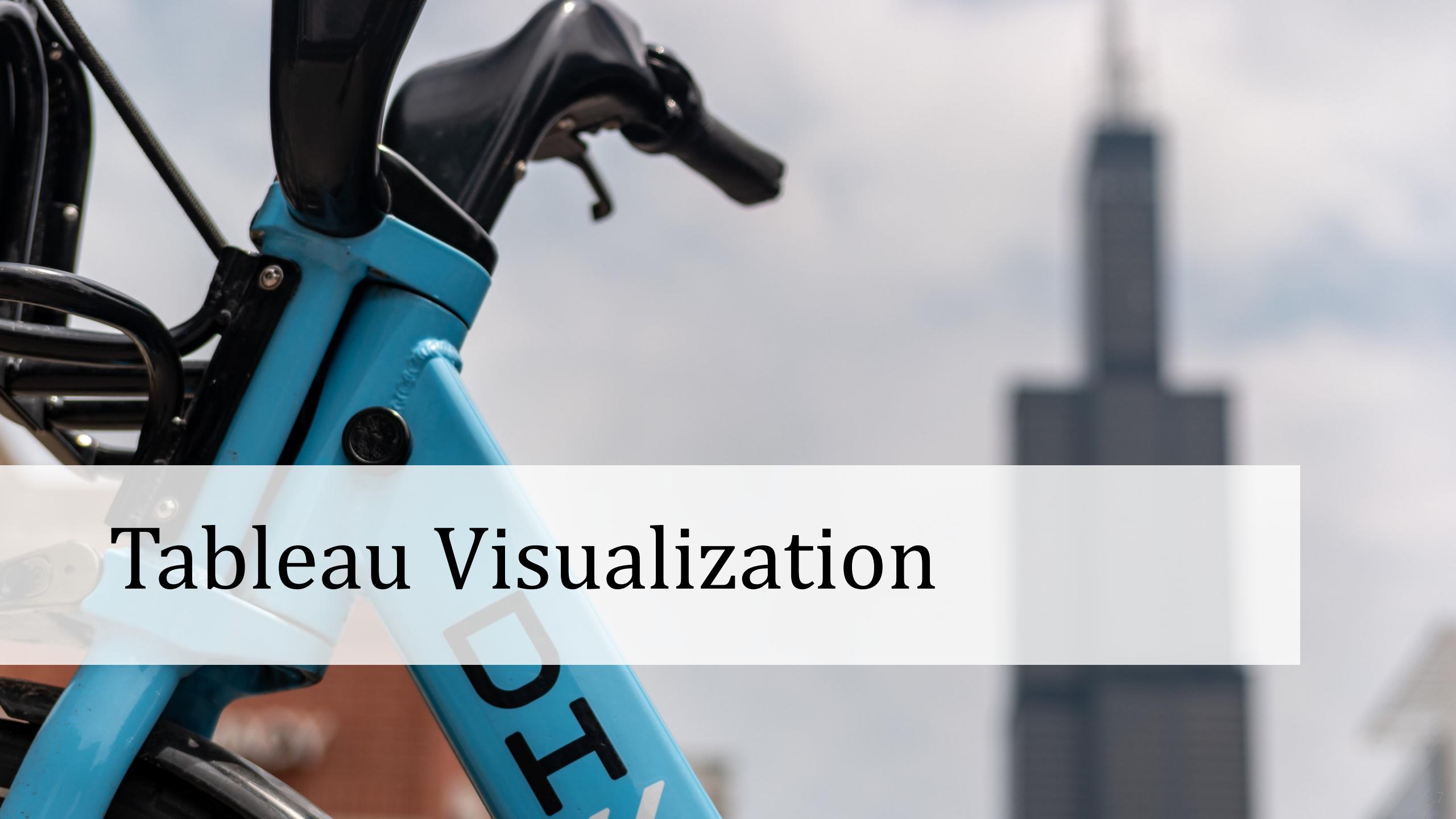
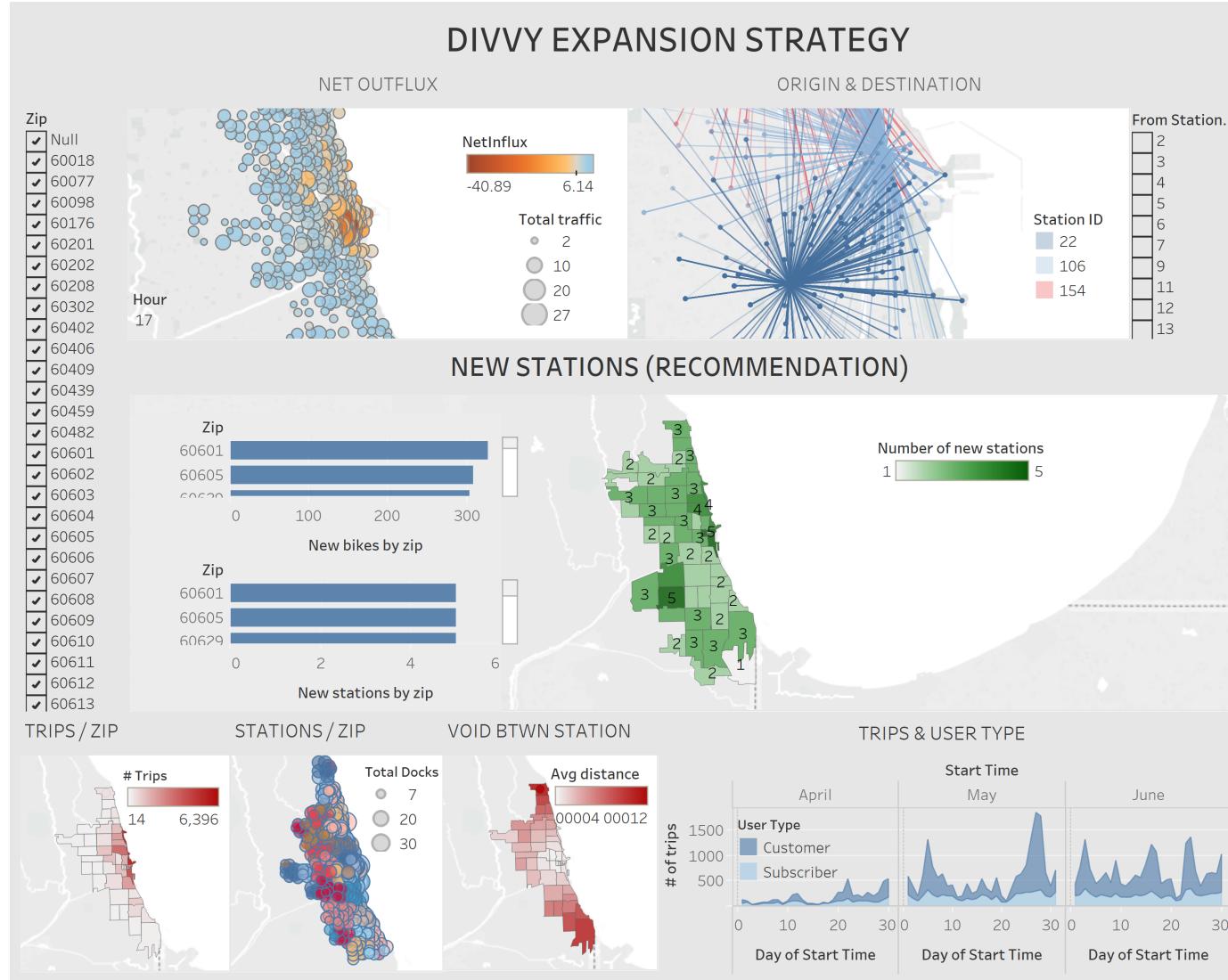


Tableau Visualization



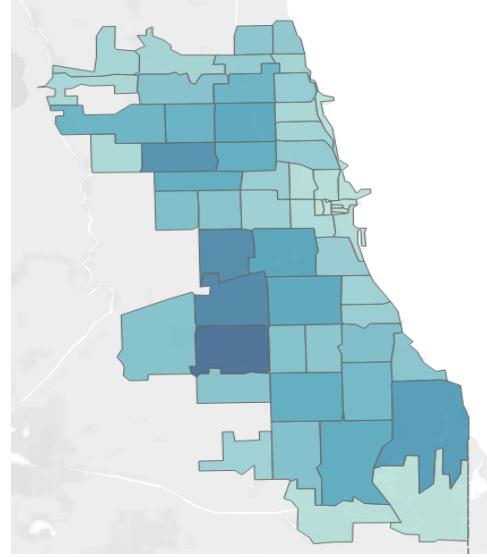
Derived recommendation from trip and zip demographics:

- **Net Outflux:** Number of bikes stalled minus number of bikes taken for each station and filtered by hour
- **Origin & Destination:** All destinations of the trips taken from a respective station
- **New Stations (Recommendation):** Suggested number of new stations per zip code, based on the previously described scoring methodology (+ Number of suggested new bikes and stations per zip code as bar chart)
- **Trips / Zip:** Average number of trips started in a respective zip code
- **Stations / Zip:** All divvy stations filtered by zip code (color wise) and number of docks (bubble size)
- **Void Btwn Station:** Average distance in 100 meters between stations within one zip code
- **Trips & User Type:** Number of trips taken filtered by subscribers and non-subscribers ('customers')

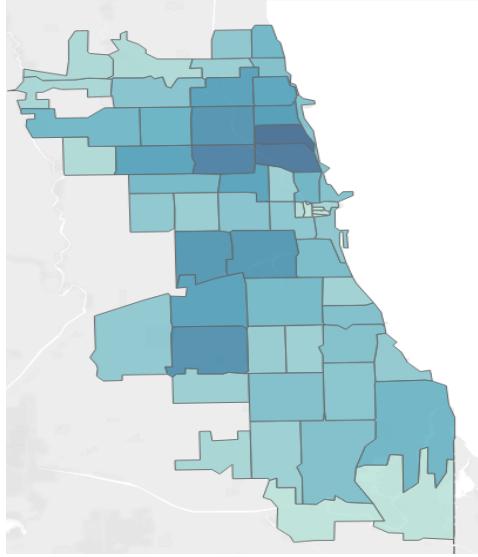
Demographics by Zip Code



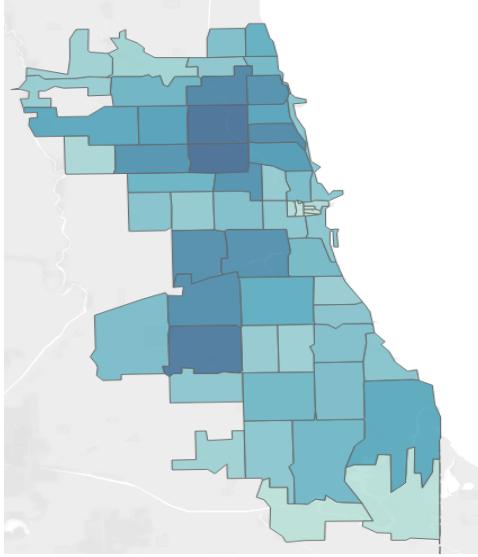
0-19 ZIP



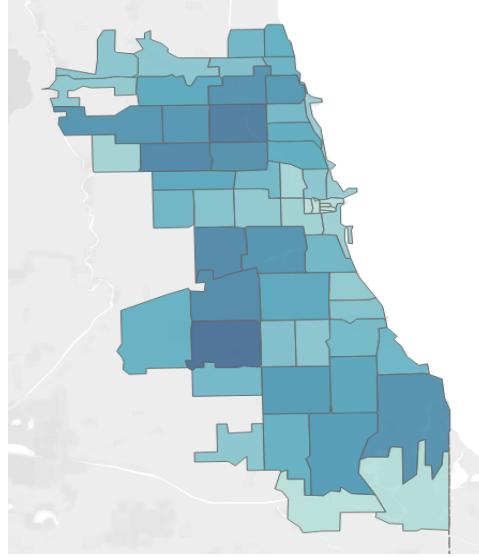
20-29



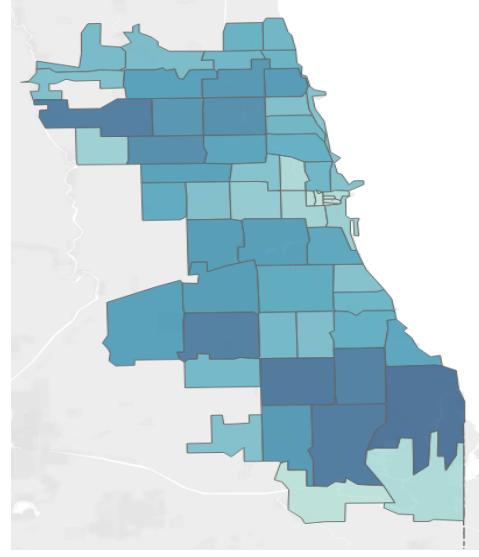
30-39



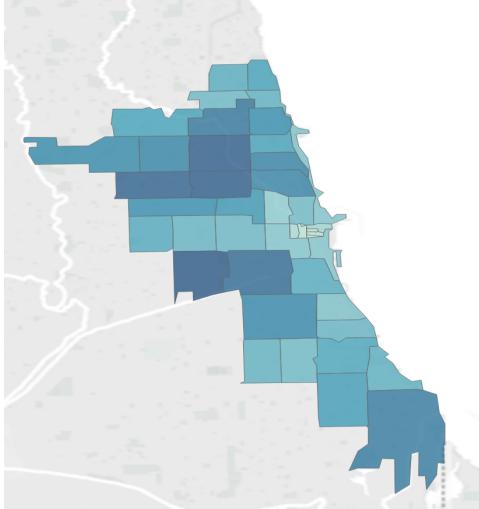
40-49



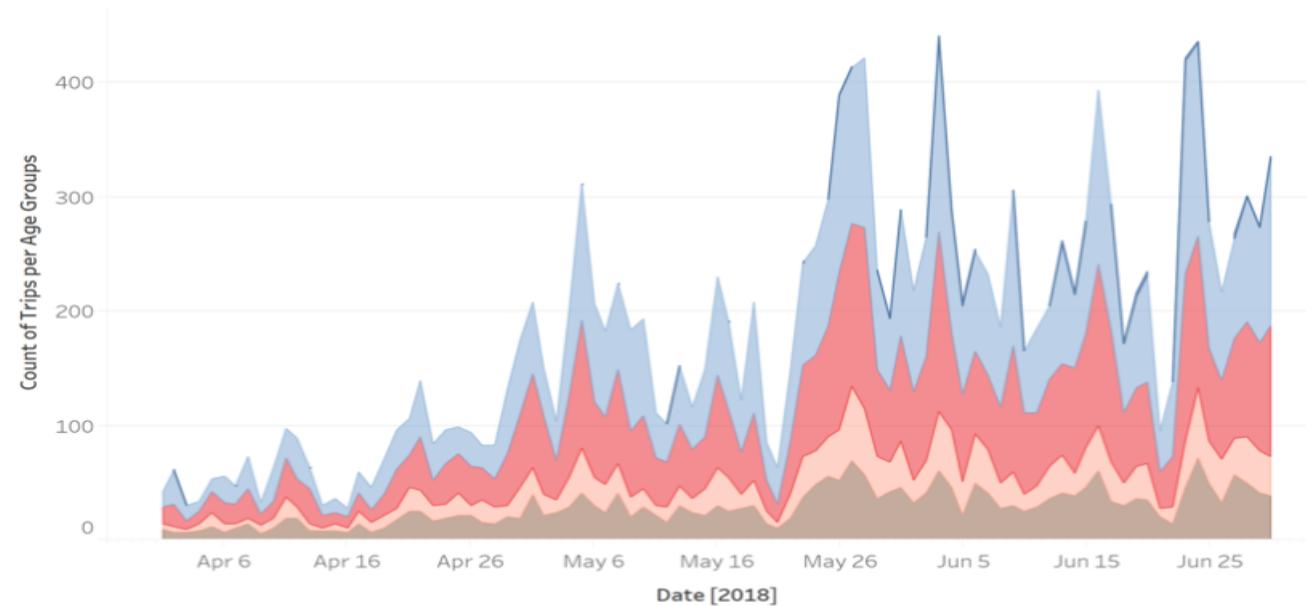
50+



Total population by zip



Number of trips taken by age groups





Summary



Recommendations and Future Vision:

- Increase stations in ZIPs farther from downtown Chicago based on scoring variables to serve the needs of local residents better
- Allocate more bikes to stations with higher net outflux (especially during summer)
- More advanced analysis based on variables like customer feedback, commercial footprints, real estate bike scores etc.
- Capitalize on the existing bike rack network in Chicago
- Expand to OLTP framework to support real time trip information.
- Scaling out to support the ever increasing data repository.

Lessons Learned:

- Choose your data sources carefully, every data source has its own conventions and business case.
- Make sure geographic data from different sources is coherent.
- Don't over normalize for OLAP - keep it simple!
- Split up data sources / use views for faster processing in tableau.
- Excel is a very powerful tool.



THANK YOU!
