



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

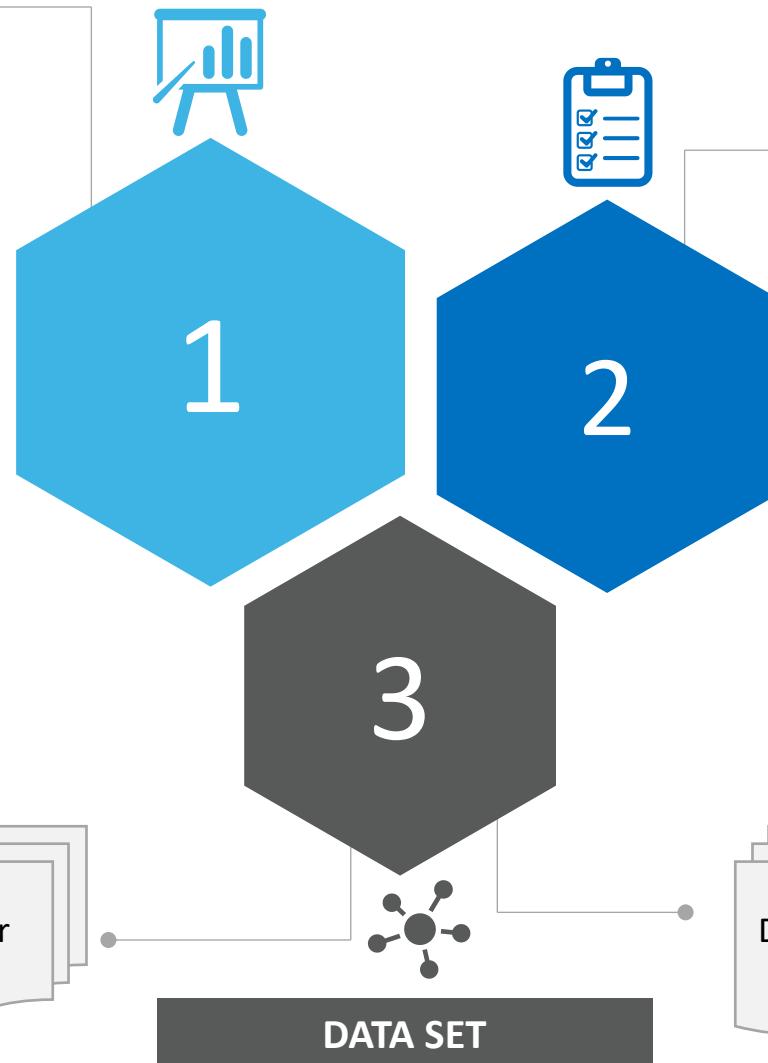
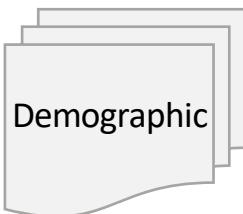
Put forward a future state blueprint for the new stations and bikes allocation process



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 (`station_id`)
) 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 (`location_id`)
) 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 (`station_id`)
) ENGINE=InnoDB;

-- Create table if not exists 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,
    `CONSTRAINT `traffic_id` FOREIGN KEY(`traffic_id`) REFERENCES `divey.dim_station`(`station_id`),
    PRIMARY KEY (`traffic_id`)
) ENGINE=InnoDB;

```

DML

```

1 * SHOW VARIABLES LIKE '%secure_file_priv%'

2 * SET @GLOBAL.secure_file_priv = '/users/tammy/Tammy/UChicago/Data_Engineering_Platform/FinalProject/Final_data/' 

3 * LOAD DATA INFILE '/users/tammy/Tammy/UChicago/Data_Engineering_Platform/FinalProject/Final_data/dim_bike_lane.csv'
4 INTO TABLE dim_bike_lane
5 FIELDS TERMINATED BY ','
6 ENCLOSED BY '\"'
7 LINES TERMINATED BY '\n'
8 IGNORE 1 ROWS;
9 
10 * SELECT * FROM dim_bike_lane;

11 * LOAD DATA INFILE '/users/tammy/Tammy/UChicago/Data_Engineering_Platform/FinalProject/Final_data/dim_bike_racks.csv'
12 INTO TABLE dim_bike_racks
13 FIELDS TERMINATED BY ','
14 ENCLOSED BY '\"'
15 LINES TERMINATED BY '\n'
16 IGNORE 1 ROWS;
17 
18 * SELECT * FROM dim_bike_racks;

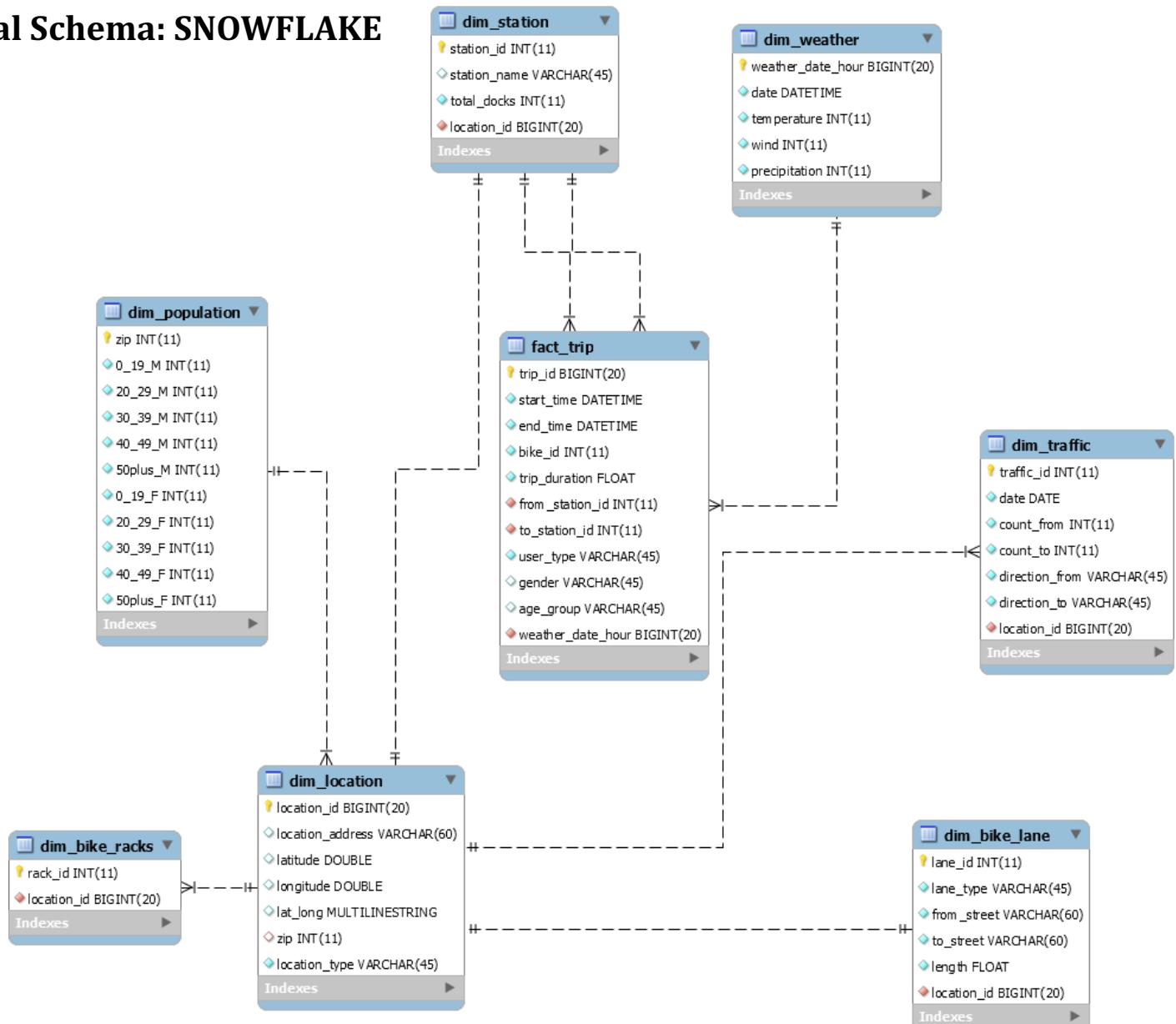
19 * LOAD DATA INFILE '/users/tammy/Tammy/UChicago/Data_Engineering_Platform/FinalProject/Final_data/dim_location_v2.csv'
20 INTO TABLE dim_location
21 FIELDS TERMINATED BY ','
22 ENCLOSED BY '\"'
23 LINES TERMINATED BY '\n'
24 IGNORE 1 ROWS;
25 
26 * SELECT * FROM dim_location;

27 * LOAD DATA INFILE '/users/tammy/Tammy/UChicago/Data_Engineering_Platform/FinalProject/Final_data/dim_weather.csv'
28 INTO TABLE dim_weather
29 FIELDS TERMINATED BY ','
30 ENCLOSED BY '\"'
31 LINES TERMINATED BY '\n'
32 IGNORE 1 ROWS;
33 
34 * SELECT * FROM dim_weather;

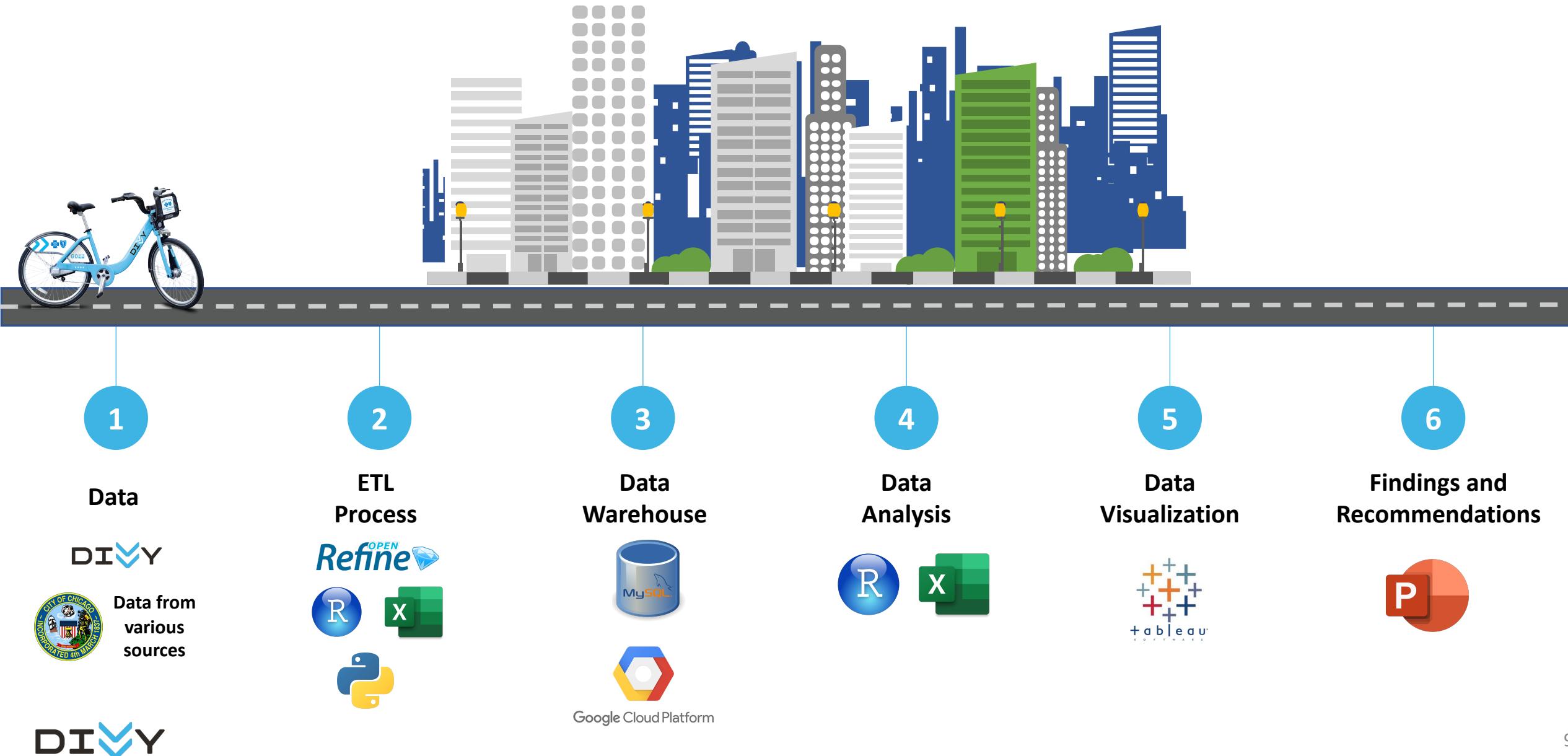
35 * LOAD DATA INFILE '/users/tammy/Tammy/UChicago/Data_Engineering_Platform/FinalProject/Final_data/fact_trip.csv'
36 INTO TABLE fact_trip
37 FIELDS TERMINATED BY ','
38 ENCLOSED BY '\"'
39 LINES TERMINATED BY '\n'
40 IGNORE 1 ROWS;
41 
42 * SELECT * FROM fact_trip;

43 * LOAD DATA INFILE '/users/tammy/Tammy/UChicago/Data_Engineering_Platform/FinalProject/Final_data/fact_trip.csv'
44 INTO TABLE fact_trip
45 FIELDS TERMINATED BY ','
46 ENCLOSED BY '\"'
47 LINES TERMINATED BY '\n'
48 IGNORE 1 ROWS;
49 
50 * SELECT * FROM fact_trip;

```



Tools



Data extraction, Cleaning, Normalization



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

```
# Number of trips by hour by weekday and weekend.
SELECT
CASE WHEN dayname(start_time) IN ("Saturday", "Sunday") THEN "Weekend" ELSE "Weekday" END AS DateType,
HOUR(start_time) AS TimeOfDay,
COUNT(trip_id) AS NoOfTrips
FROM fact_trip
GROUP BY TimeOfDay, DateType
ORDER BY TimeOfDay DESC;
```

```
SELECT
# #. Number of TripIn per zip
    tripft.station_id AS stationId,
    tripft.start_time AS TimeOfDay,
    tripft.end_time AS TimeOfDay,
    tripft.tripTo AS tripFrom,
    tripft.tripTo,
    tripft.tripTo - (tripft.tripTo - tripft.tripFrom) AS NetTrips,
    (tripft.totalBikes - tripft.tripFrom + tripft.tripTo) AS docksNeededDpss
FROM
    fact_trip ft
    JOIN dim_station ds ON ft.station_id = ds.station_id
    LEFT JOIN dim_location dl ON dl.location_id = ds.location_id
GROUP BY zip
ORDER BY TripIn DESC;
```

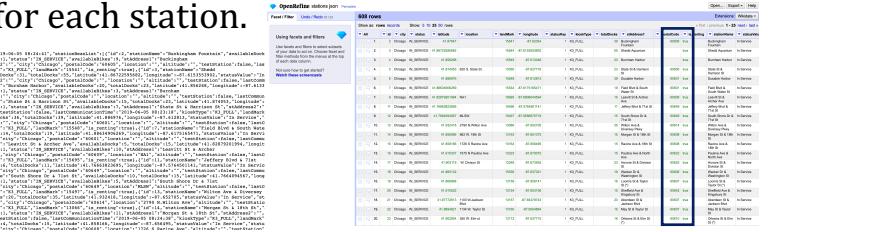


- Clean all dimensional tables to import to MySQL
- Analyze descriptive data: customer profiling, zip, stations
- Build the scoring system for research objectives' purpose: add more stations and bikes.

community_id	community_area	Male	0-19		20-29		30-39		40-49		50-59		60-69		70-79		80-89		90-99					
			Female	Total	Traffic	Bike Parks	Dock Stations																	
1	Rogers Park	27752	27235	12334	100000	45700	17204	8000	500	15	47	60,80	10,80	90,75	50,75	92,7	67	20,81	38	31,74	(1), 31	114	26,71	24,71
2	Wicker Park	35322	36400	20005	100000	43000	16300	4500	300	15	48	60,80	10,80	90,75	50,75	91,73	63	14,83	38	31,74	51,74	10,83	14,83	21,74
3	Uptown	29512	26850	8106	100000	42000	15700	4000	300	15	49	60,80	10,80	90,75	50,75	90,74	62	15,83	38	31,74	51,74	10,83	15,83	21,74
4	Lincoln Square	19309	20184	7168	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	91,75	53	15,83	38	31,74	51,74	10,83	15,83	21,74
5	North Center	15557	16310	6289	100000	42000	15700	4000	300	15	47	60,80	10,80	90,75	50,75	90,76	52	15,84	38	31,74	51,74	10,84	15,84	21,74
6	Lake View	47000	47368	9334	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,77	53	15,85	38	31,74	51,74	10,85	15,85	21,74
7	Lincoln Park	30430	33686	10135	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,78	53	15,86	38	31,74	51,74	10,86	15,86	21,74
8	New North Side	37337	43147	6842	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,79	53	15,87	38	31,74	51,74	10,87	15,87	21,74
9	Edison Park	5307	5828	2556	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,80	53	15,88	38	31,74	51,74	10,88	15,88	21,74
10	North Park	23222	32101	5327	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,81	53	15,89	38	31,74	51,74	10,89	15,89	21,74
11	Jefferson Park	12393	13155	5707	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,82	53	15,90	38	31,74	51,74	10,90	15,90	21,74
12	Forest Glen	8934	9574	4850	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,83	53	15,91	38	31,74	51,74	10,91	15,91	21,74
13	North Park	8649	9282	4429	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,84	53	15,92	38	31,74	51,74	10,92	15,92	21,74
14	Albany Park	26407	25153	14526	100000	42000	15700	4000	300	15	47	60,80	10,80	90,75	50,75	90,85	53	15,93	38	31,74	51,74	10,93	15,93	21,74
15	Portage Park	31337	32787	16227	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,86	53	15,94	38	31,74	51,74	10,94	15,94	21,74
16	Irving Park	26674	26686	13426	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,87	53	15,95	38	31,74	51,74	10,95	15,95	21,74
17	Dunning	20504	21428	5943	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,88	53	15,96	38	31,74	51,74	10,96	15,96	21,74
18	Morgan Park	6032	6838	3452	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,89	53	15,97	38	31,74	51,74	10,97	15,97	21,74
19	Belmont-Cragin	39609	38384	25998	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,90	53	15,98	38	31,74	51,74	10,98	15,98	21,74
20	Hermosa	12566	12444	8416	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,91	53	15,99	38	31,74	51,74	10,99	15,99	21,74
21	Avondale	20119	19143	10710	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,92	53	15,100	38	31,74	51,74	10,100	15,100	21,74
22	Logan Square	36805	35986	16711	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,93	53	15,101	38	31,74	51,74	10,101	15,101	21,74
23	Humboldt Park	27733	29505	19317	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,94	53	15,102	38	31,74	51,74	10,102	15,102	21,74
24	West Town	42058	40178	13909	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,95	53	15,103	38	31,74	51,74	10,103	15,103	21,74
25	Austin	45189	53325	30332	100000	42000	15700	4000	300	15	48	60,80	10,80	90,75	50,75	90,96	53	15,104	38	31,74	51,74	10,104	15,104	21,74
26	



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



- 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...



- 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 mySQL.



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

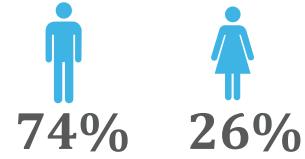


Users Type

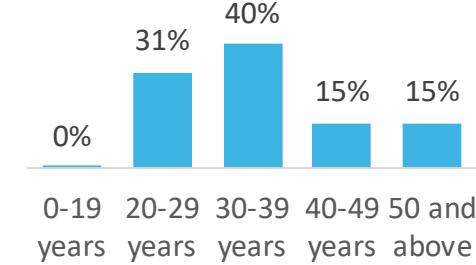


■ Subscriber
■ Non-subscriber

Gender



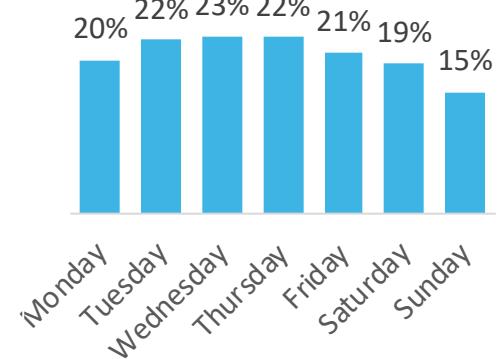
Age Group



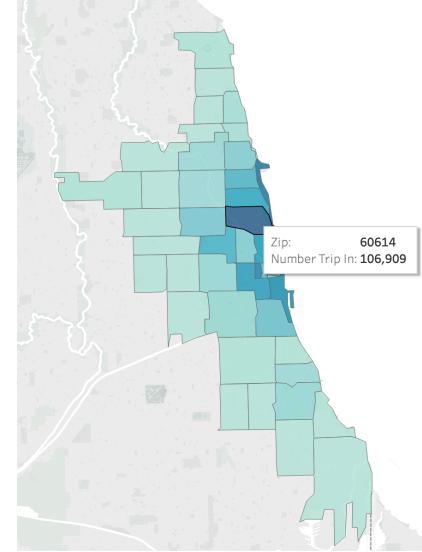
Average Distance Travel



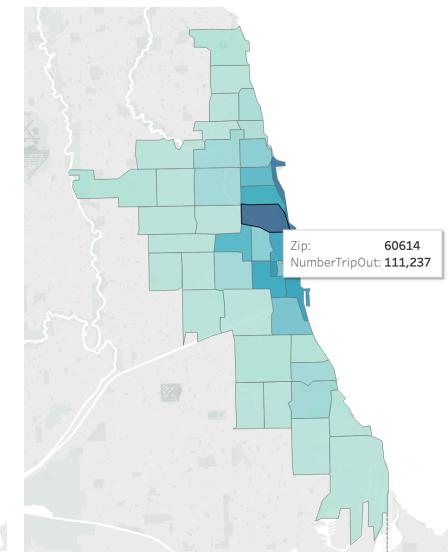
Trip by day



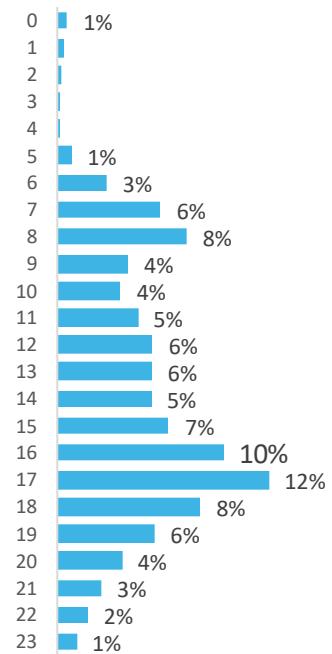
Trip-in by area



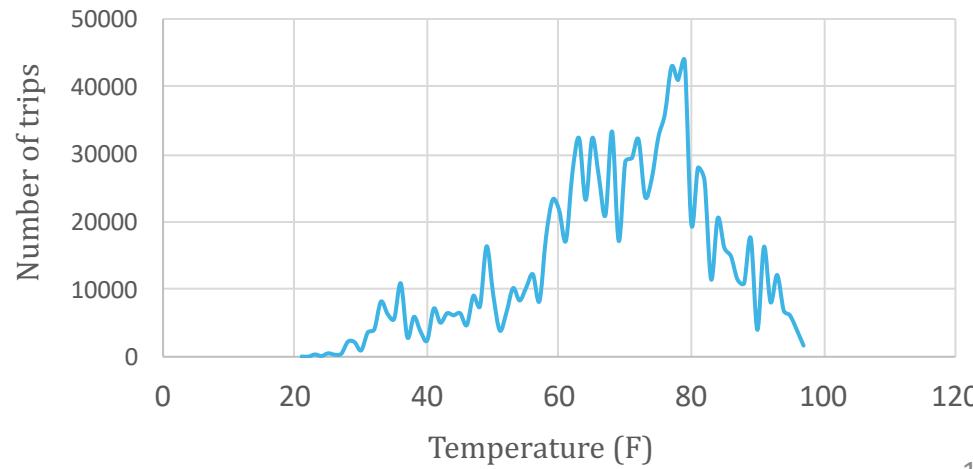
Trip-out by area



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			% of total trips	Subscriber %										
		Age					Number of stations	Total # of docks	Avg # of docks	Avg of Avg Distance from other stations (miles)	Trips Out	Trips In	Net											
		Male%	Female%																					
60605	1115	48.5%	51.5%	0.0678	0.138	0.0294	0.0585	0.0134	0.0719	0.1451	0.1683	0.0603	0.206	8,300	356	28.9	68,302	65,243	(3.65%)	6.33%	60.7%			
60601	1115	48.4%	50.6%	0.0659	0.1702	0.0933	0.0582	0.0103	0.0716	0.1744	0.1693	0.0593	0.173	23,800	1	33.1	4,59	68,506	63,594	(4.91%)	6.23%	72.4%		
60609	6436	50.1%	49.8%	0.0651	0.1272	0.1173	0.0563	0.0549	0.0634	0.1529	0.0785	0.0723	0.0769	0.113	8,75	28	34.4	6,25	2,242	2,242	0.4%	61.9%	62.0%	
60649	13,664	52.4%	47.5%	0.0651	0.1621	0.0696	0.0527	0.0951	0.0627	0.181	0.0375	0.051	0.043	21,700	195	34	638	18.8	4,82	106,936	111,237	4,328	10.23%	73.3%
60614	66623	47.6%	52.4%	0.0571	0.1642	0.0696	0.0527	0.0951	0.0627	0.181	0.0375	0.051	0.043	21,700	195	27	370	13.7	5.03	14,489	15,288	759	1.4%	67.8%
60608	82743	47.6%	52.4%	0.1505	0.1084	0.0693	0.0569	0.0367	0.1355	0.0968	0.0916	0.0503	0.104	18,600	109	45	495	18.6	4,57	45,14	46,260	146	4.3%	66.0%
60622	53533	51.5%	48.3%	0.0892	0.1414	0.0354	0.0593	0.0666	0.0845	0.1352	0.1295	0.0566	0.031	34,400	354	25	102	1.14	4,71	37,347	38,245	930	2.6%	59.7%
60609	77,034	51.5%	48.3%	0.0892	0.1414	0.0354	0.0593	0.0666	0.0845	0.1352	0.1295	0.0566	0.031	34,400	354	1	1	1.14	4,71	37,347	38,245	930	2.6%	59.7%
60607	23,850	50.8%	49.1%	0.0871	0.1485	0.1337	0.0522	0.0206	0.0893	0.1536	0.178	0.054	0.036	29,300	62	26	453	17.4	4,40	61,385	61,617	22	3.8%	31.1%
60642	18465	51.3%	48.3%	0.0858	0.1424	0.1361	0.0503	0.1303	0.0568	0.0804	0.1100	0.9	12	217	18.1	4,37	24,446	24,523	107	2.3%	90.8%			
60610	37730	48.4%	51.5%	0.0897	0.1339	0.0951	0.0553	0.0456	0.1057	0.067	0.0892	0.22	100	122	18	403	14.2	4,52	61,694	62,402	708	5.8%	81.3%	
60624	28722	48.4%	51.5%	0.0834	0.1321	0.0975	0.0533	0.0446	0.1058	0.0659	0.0887	0.18	100	149	16	23	23.3	4,67	83,207	95,133	9,326	7.7%	58.4%	
60604	47,934	48.5%	51.5%	0.0704	0.1056	0.0584	0.0584	0.0738	0.0587	0.0522	0.0522	0.113	0.01A	11,000	15	14	558	10.6	4,62	15,230	15,230	4,834	33.2%	58.4%
60603	497	49.3%	50.7%	0.0704	0.163	0.0526	0.0592	0.0707	0.0724	0.1751	0.0946	0.0522	0.127	13,700	98	5	135	2.0	4,56	27,542	24,721	(2.82%)	2.47%	60.7%
60616	24,791	51.3%	48.6%	0.0598	0.0954	0.0776	0.0622	0.0499	0.1025	0.036	0.0588	0.0684	0.026	6,100	87	29	447	15.4	5.35	32,058	32,470	323	3.0%	60.7%
60602	7,734	48.3%	51.6%	0.0622	0.1666	0.1057	0.0524	0.0605	0.0913	0.0562	0.0542	0.0589	0.032	12,700	122	12	354	15.5	4.34	78,911	78,911	323	1.6%	60.7%
60637	49,008	48.4%	51.5%	0.1343	0.1704	0.0586	0.0713	0.1738	0.0587	0.0522	0.0522	0.113	0.01A	20,900	72	17	264	15.5	8.15	12,584	12,429	(259)	1.1%	81.1%
60657	66601	49.7%	50.3%	0.0525	0.1662	0.1753	0.0585	0.0675	0.0532	0.166	0.165	0.0565	0.0865	12,900	201	20	371	16.6	5.28	54,456	57,328	2,872	5.2%	80.3%
60647	82,937	50.4%	49.6%	0.0516	0.1672	0.1932	0.0513	0.0709	0.0587	0.1603	0.0642	0.0672	0.054	10,600	200	19	333	16.6	5.43	26,516	27,351	835	2.4%	82.6%
60621	25,200	51.3%	48.6%	0.0652	0.1466	0.1239	0.0536	0.0553	0.0739	0.0583	0.0583	0.0583	0.026	15,800	64	11	177	16.1	7.54	10,823	10,307	94	1.0%	76.7%
60615	40,008	49.3%	50.6%	0.0591	0.1374	0.0653	0.0523	0.0536	0.0739	0.0583	0.0583	0.0583	0.026	18,000	111	32	183	16.1	7.54	10,823	10,307	94	1.0%	76.7%
60618	32,085	50.2%	49.8%	0.1238	0.1957	0.107	0.0703	0.107	0.1223	0.0595	0.0394	0.0697	0.0797	18,700	193	21	308	14.7	5.95	10,373	11,778	806	1.0%	85.2%
60613	48,935	50.4%	49.6%	0.0502	0.166	0.176	0.056	0.0626	0.0708	0.0593	0.0593	0.0593	0.026	11,600	76	23	426	18.5	5.90	43,431	44,295	864	4.1%	79.5%
60623	23,912	50.5%	49.5%	0.0502	0.166	0.176	0.056	0.0626	0.0708	0.0593	0.0593	0.0593	0.026	13,700	73	13	123	12.3	6.62	11,101	12,217	176	7.3%	80.8%
60625	42,052	50.5%	49.5%	0.0502	0.166	0.176	0.056	0.0626	0.0708	0.0593	0.0593	0.0593	0.026	13,700	73	13	123	12.3	6.62	11,101	12,217	176	7.3%	80.8%
60640	65,736	51.8%	48.2%	0.0775	0.1193	0.0783	0.0399	0.0727	0.1031	0.0332	0.0722	0.1301	0.026	16,700	205	14	282	20.1	6.71	22,613	23,060	447	2.1%	80.4%
60619	63,830	48.5%	51.4%	0.1177	0.0525	0.0595	0.0572	0.0674	0.0642	0.074	0.1203	17,800	31	17	183	10.8	9.86	638	684	46	0.06%	77.6%		
60634	74,302	49.3%	50.3%	0.1233	0.0657	0.0711	0.0589	0.0893	0.1273	0.0721	0.0736	0.0708	0.026	17,800	112	1	1	15.0	5.33	3,164	3,653	523	0.32%	73.8%
60620	53,242	50.2%	49.6%	0.0517	0.0654	0.0705	0.0543	0.0717	0.0604	0.0623	0.0673	0.0745	0.026	15,200	28	13	130	8.18	8.28	1,591	2,311	731	1.0%	88.3%
60621	35,918	50.8%	49.1%	0.0517	0.0654	0.0705	0.0543	0.0717	0.0604	0.0623	0.0673	0.0745	0.026	15,200	28	13	130	8.18	8.28	1,591	2,311	731	1.0%	88.3%
60624	38,303	48.6%	51.3%	0.051	0.0653	0.0607	0.0589	0.0613	0.1844	0.0762	0.0694	0.0674	0.0534	11,800	35	8	89	11.0	5.71	3,330	320	(19)	0.03%	68.1%
60623	32,112	53.7%	46.4%	0.1844	0.0991	0.0601	0.0694	0.0732	0.0644	0.0674	0.0534	0.0534	0.026	13,700	109	6	66	11.0	5.88	560	588	28	0.05%	82.0%
60645	45,288	50.5%	49.5%	0.1342	0.1703	0.0763	0.0689	0.0404	0.1369	0.0773	0.0774	0.0698	0.0534	11,400	55	7	65	12.1	8.99	1,236	1,297	611	7.2%	80.2%
60639	44,593	50.5%	49.5%	0.1342	0.1703	0.0763	0.0689	0.0404	0.1369	0.0773	0.0774	0.0698	0.0534	11,400	55	4	44	7.93	7.93	5,057	5,057	0	0.05%	82.0%
60660	47,521	51.2%	48.8%	0.0844	0.1033	0.0563	0.0791	0.0806	0.0982	0.0512	0.0753	0.0324	0.026	35,200	42	5	31	18.2	7.57	5,215	5,686	486	1.01%	79.7%
60641	7,868	49.8%	50.4%	0.0791	0.0803	0.0687	0.0612	0.0737	0.0802	0.0875	0.0708	0.0738	0.026	39,900	81	6	78	13.0	7.21	1,173	1,110	(63)	0.1%	77.2%
60630	54,093	49.2%	50.3%	0.1221	0.0795	0.0768	0.0757	0.0899	0.1239	0.0733	0.0789	0.0741	0.0653	12,300	37	2	26	19.1	7.28	230	239	9	0.02%	73.5%
60651	49,552	49.7%	50.2%	0.0794	0.0858	0.0658	0.0624	0.0744	0.0814	0.0861	0.0784	0.0784	0.026	16,200	37	6	65	20.5	6.25	2,942	2,795	147	0.02%	69.3%
60644	49,552	49.7%	50.2%	0.0794	0.0858	0.0658	0.0624	0.0744	0.0814	0.0861	0.0784	0.0784	0.026	16,200	37	13	134	10.7	6.25	2,942	2,795	147	0.02%	69.3%
60636	49,023	48.6%	51.3%	0.1537	0.0871	0.0527	0.0587	0.0818	0.0785	0.0603	0.0873	0.0573	0.026	18,200	20	8	80	10.0	8.10	105	109	4	0.01%	76.2%
60617	84161	48.5%	51.5%	0.0513	0.0678	0.0563	0.0643	0.0681	0.0693	0.0681	0.0681	0.0681	0.026	8,600	190	6	62</td							

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% -

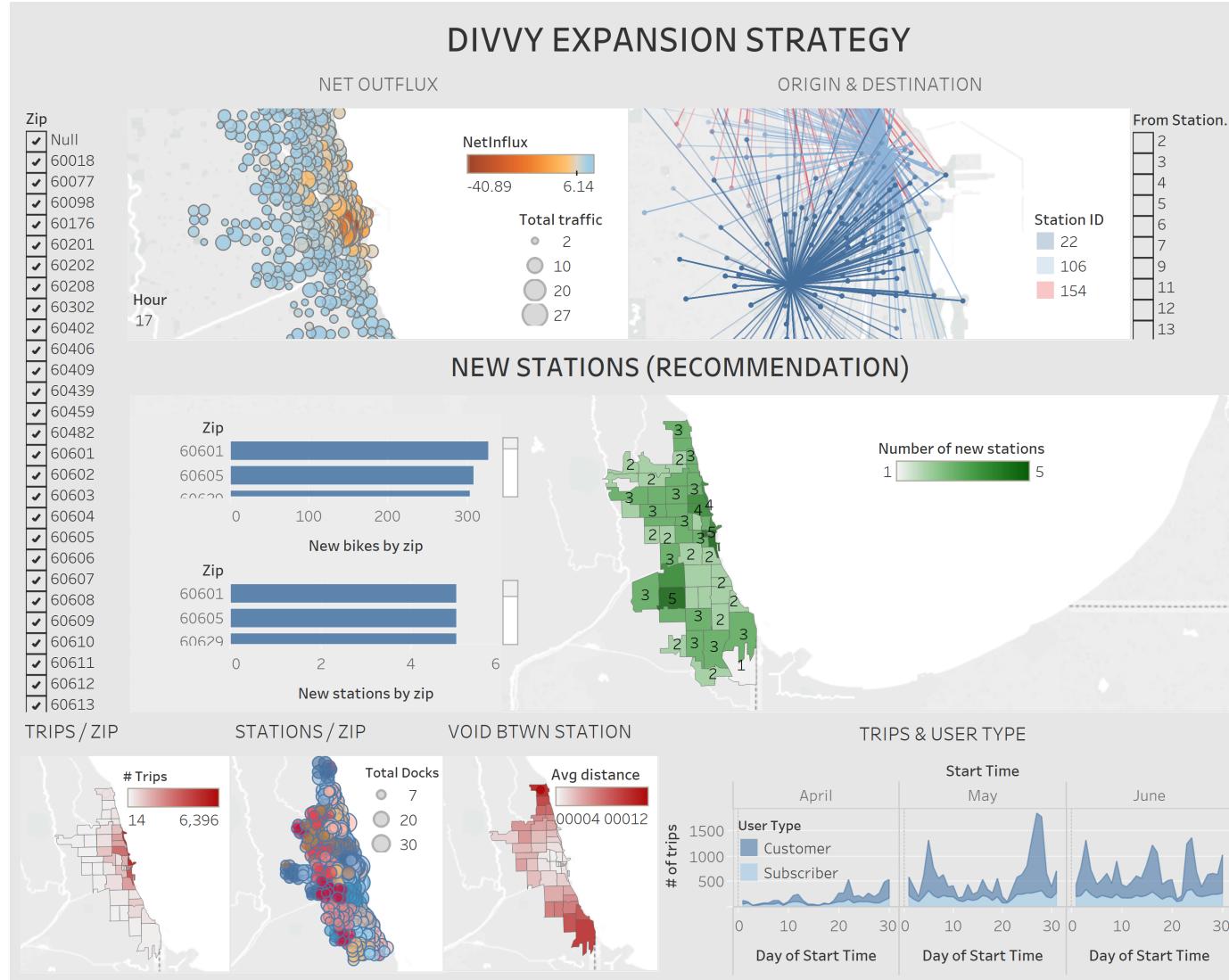
Scores by zip code



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%	20_39	Vehicle Volume	Number of stations	Area of avg distance from other stations (mi)	Trips	Net (Trips Out)	10%	5%	10%	5%	20%	10%	5%	10%	5%	10%	5%	10%	5%		
	Total	Male%	20_39	Vehicle Volume	Number of stations	Area of avg distance from other stations (mi)	Trips	Net (Trips Out)	10%	5%	10%	5%	20%	10%	5%	10%	5%	10%	5%	10%	5%		
Weight	15%	5%	10%	5%	5%	5%	10%	10%	10%	5%	10%	5%	10%	10%	10%	10%	5%	10%	5%	10%	5%		
Total Points	1500	500	1000	500	500	2000	1000	1000	1000	500	1000	500	1000	1000	1000	1000	500	1000	500	1000	500		
60605	15072	48.54	0.9394	8,300	356	19	4.75	68,902	(3,659)	60.7%	12.7	8.7	23.9	3.9	(3,26)	15.4	65.0	84.3	24.4	240.0	5	308	
60601	1115	43.45	0.5329	23,800	1	11	4.59	88,506	(4,912)	72.4%	6.2	8.8	25.0	11.1	0.1	(8,09)	14.3	64.7	103.2	23.2	255.0	5	328
60609	14591	50.25	0.3668	8,100	75	28	6.49	29,953	145	87.5%	36.1	9.0	14.4	3.8	6.7	(46,05)	21.1	2.4	-	35.2	62.6	2	106
60649	46864	43.85	0.2693	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
60614	66623	47.63	0.313	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
60608	82743	52.45	0.3789	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
60622	62565	61.5	0.5416	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
60606	2314	48.25	0.5328	10,200	227	6	4.40	37,508	(3,42)	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
60607	23902	49.25	0.5333	23,300	62	26	4.40	61,385	232	91.6%	13.3	8.8	25.0	13.7	5.6	(42,76)	14.3	57.3	-	36.7	132.5	3	170
60642	18485	51.5	0.5452	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
60610	37730	46.45	0.4958	22,100	122	18	4.52	16,594	708	81.9%	210	8.3	22.8	10.3	10.9	(29,61)	14.7	58.2	-	33.0	149.6	3	192
60611	28722	48.45	0.4955	18,100	149	16	4.67	83,207	5,326	58.4%	16.0	8.3	22.8	8.5	13.3	(26,32)	15.2	84.2	-	23.5	165.4	4	213
60654	14880	46.75	0.4893	23,000	115	14	4.40	68,187	(2,80)	88.8%	8.3	8.3	22.9	10.8	10.3	(23,03)	14.3	64.4	64.6	35.6	216.7	5	278
60604	575	49.45	0.5338	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
60603	497	49.35	0.513	13,700	38	5	4.56	27,542	(2,82)	55.5%	0.3	8.8	24.9	6.4	8.8	(2,22)	14.8	26.0	65.0	22.4	163.2	4	217
60616	48437	48.15	0.3694	6,100	87	29	5.35	32,058	412	77.6%	27.0	8.6	16.3	7.8	7.8	(47,70)	17.4	30.3	-	31.2	94.3	2	121
60602	1210	49.35	0.5322	18,750	50	3	4.53	17,152	(32)	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
60651	7298	48.25	0.5311	15,600	102	12	4.34	78,847	(3,43)	92.2%	4.3	8.8	25.0	7.3	9.1	(9,74)	14.1	74.4	56.1	37.6	217.0	5	279
60637	49008	44.95	0.3076	20,900	72	17	8.15	12,684	(3,55)	81.4%	27.5	8.0	14.4	8.8	6.4	(27,89)	26.5	12.0	5.9	32.9	115.4	2	148
60657	69001	49.75	0.696	12,900	201	20	5.28	54,456	2,372	80.8%	36.7	8.3	28.4	8.0	18.0	(22,89)	17.1	51.4	-	32.4	166.2	4	214
60647	82737	50.45	0.4962	16,600	204	24	5.43	26,516	835	82.6%	48.6	9.0	21.6	5.0	18.3	(39,47)	17.6	25.0	-	33.3	138.9	3	178
60612	32478	48.45	0.4944	29,500	76	19	4.68	13,200	(146)	91.5%	18.6	8.7	20.7	13.8	6.8	(31,25)	15.2	12.5	14.9	36.8	116.7	2	150
60615	40608	44.95	0.3691	10,800	64	11	7.94	10,823	84	76.7%	22.6	8.0	16.9	5.1	5.7	(18,09)	24.5	10.2	-	30.9	105.6	2	136
60618	92089	50.25	0.3387	18,700	193	21	5.95	10,973	805	85.2%	9.0	8.5	18.5	8.7	17.3	(34,54)	18.3	10.4	-	34.3	134.2	3	172
60613	42825	50.45	0.556	16,600	76	23	5.90	43,431	864	89.8%	26.9	9.0	26.1	5.4	6.8	(27,83)	18.2	41.0	-	32.0	128.6	3	165
60653	23912	43.25	0.5031	19,900	59	9	6.62	2,101	176	73.3%	16.5	7.7	14.2	9.3	5.3	(14,80)	21.5	2.0	-	23.5	91.4	2	117
60625	78654	50.05	0.3928	23,800	286	19	6.85	11,857	390	87.8%	4.38	8.3	18.4	11.1	25.6	(31,25)	22.2	11.2	-	35.4	145.5	3	187
60640	65736	51.85	0.4281	16,700	205	14	6.71	22,613	447	80.4%	36.6	9.3	20.1	7.8	18.4	(23,03)	21.8	21.3	-	32.4	144.6	3	186
60619	63830	43.85	0.2334	17,800	31	17	9.86	638	46	77.6%	35.5	7.8	11.0	8.3	2.8	(27,86)	32.0	0.6	-	31.2	101.3	2	130
60634	74302	49.15	0.2865	21,800	112	1	5.33	3,164	529	79.8%	8.8	8.8	13.4	10.2	10.0	(164)	17.3	3.0	-	32.1	134.6	3	173
60626	50144	50.45	0.4019	7,100	145	15	8.92	7,235	60	84.5%	27.9	9.0	18.9	3.3	13.0	(24,67)	29.0	6.8	-	34.0	117.3	3	151
60621	35916	44.85	0.2988	15,000	28	12	8.28	395	(14)	86.2%	20.0	8.0	12.1	7.0	2.5	(19,74)	26.8	0.4	0.3	34.8	92.3	2	119
60624	38103	46.85	0.2786	11,800	35	8	5.71	339	(19)	68.1%	21.2	8.3	13.1	5.5	3.1	(3,85)	18.5	0.3	0.4	27.4	84.8	2	109
60623	3212	53.75	0.3391	13,700	109	6	5.88	560	28	82.0%	51.3	8.6	15.9	6.4	9.8	(3,87)	18.1	0.5	-	33.0	135.7	3	174
60645	45280	49.85	0.3074	11,400	55	7	8.93	1296	2	77.2%	25.2	8.3	14.4	5.3	4.3	(11,51)	29.2	12	-	31.1	108.7	2	140
60659	30109	43.15	0.2895	44,000	19	4	7.79	601	96	75.5%	212	8.8	13.5	20.6	17	(5,58)	25.3	0.6	-	30.4	115.5	2	148
60650	51257	51.25	0.3896	35,200	42	5	7.57	5,215	471	79.7%	238	9.2	18.2	16.5	3.8	(2,22)	24.6	4.9	-	32.1	124.8	3	160
60641	78659	48.85	0.2311	28,900	91	6	7.21	1,173	(52)	77.2%	39.9	8.3	15.1	18.7	7.3	(9,87)	23.4	11	15	31.1	137.0	3	176
60630	54099	43.25	0.3003	12,900	37	2	7.28	230	9	73.8%	30.1	8.8	14.1	5.8	3.3	(2,29)	23.8	0.2	-	32.6	112.2	2	144
60651	84272	47.75	0.2911	24,900	31	4	6.25	804	(7)	80.6%	35.8	8.5	13.7	16.3	2.8	(6,58)	20.3	0.8	0.2	32.2	119.9	3	154
60644	48652	45.95	0.2361	3,900	49	13	6.89	278	(59)	69.4%	21.1	8.2	12.5	4.6	4.4	(1,39)	22.4	0.3	0.4	28.6	86.4	2	111
60636	40923	46.85	0.257	18,200	20	8	8.10	105	4	76.2%	22.8	8.3	12.1	8.5	1.8	(1,16)	18.6	0.1	-	30.7	97.4	2	125
60617	46.85	42.504	8,600	190	6	10.84	191	(2)	90.1%	8.3	11.8	4.0	14.7	5.8	4.0	(1,64)	18.6	0.2	0.2	30.4	132.5	3	170
60201	#NA	#NA	#NA	#NA	9	1162	3,783	(

Tableau Visualization



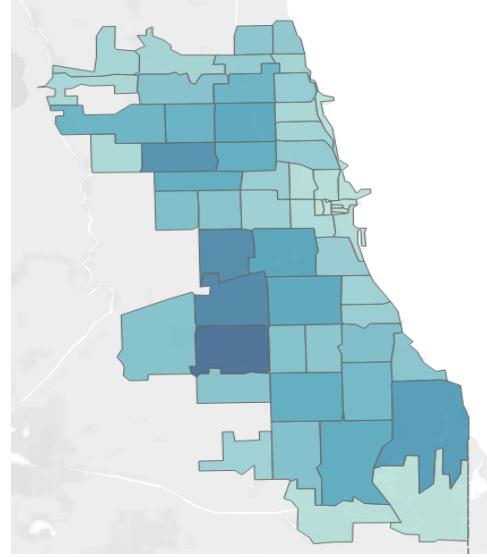
Derived recommendation from trip and zip demographics:

- **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')

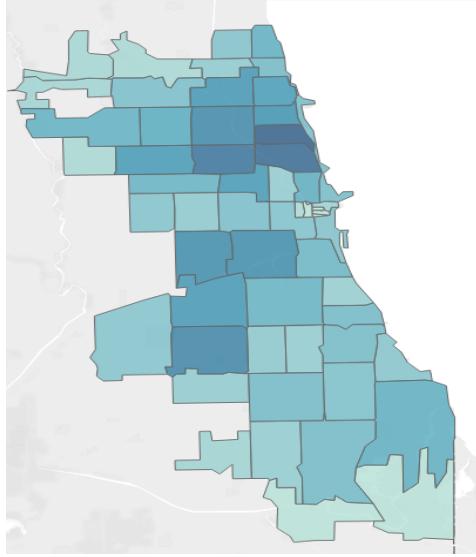
Demographics by Zip Code



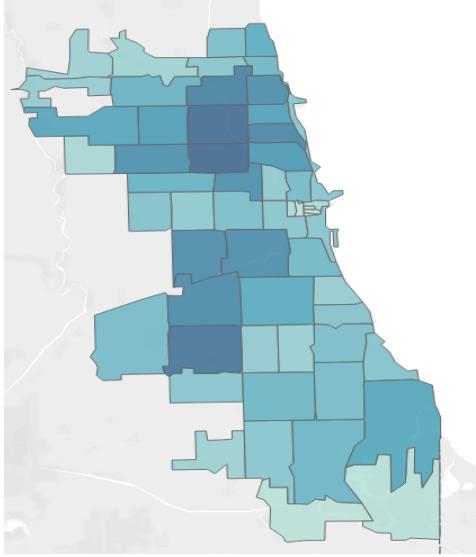
0-19 ZIP



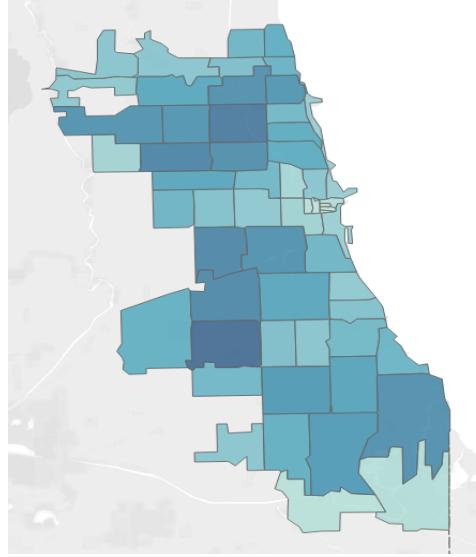
20-29



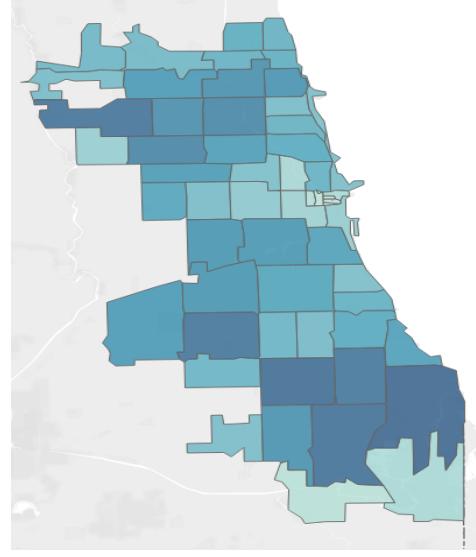
30-39



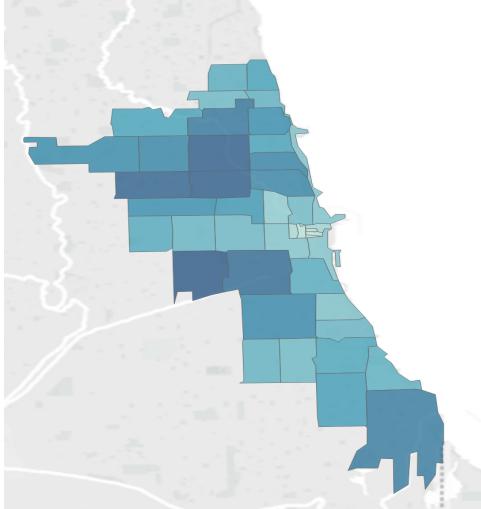
40-49



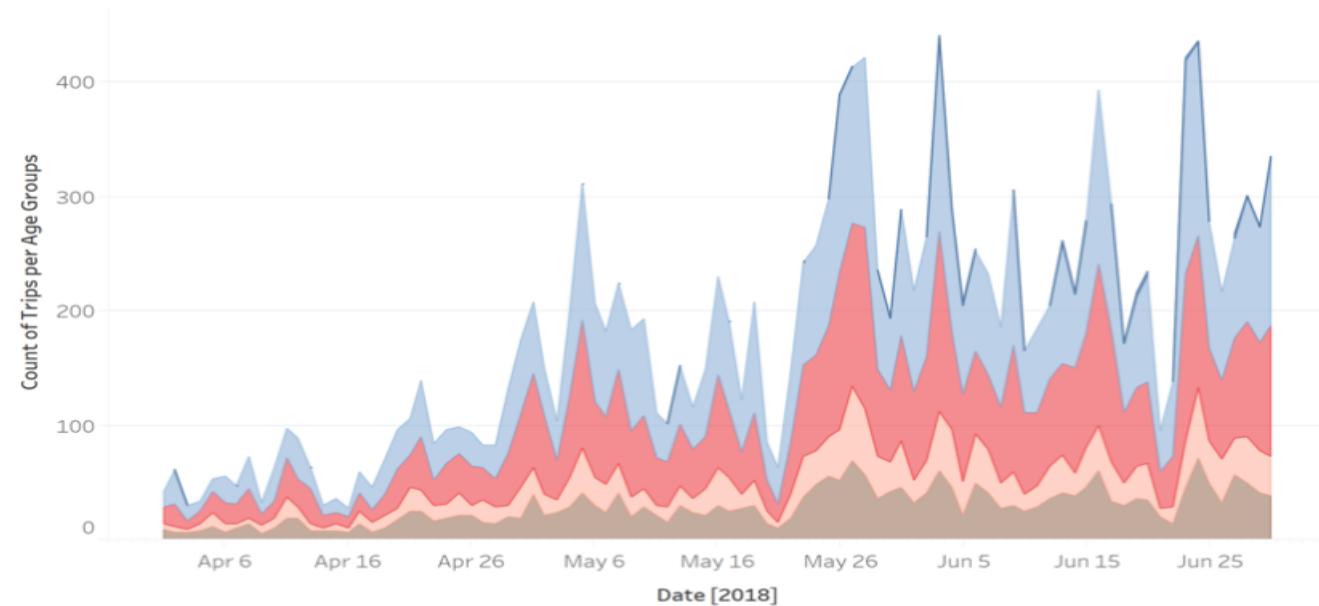
50+



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!
