



DI  **Y**

**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



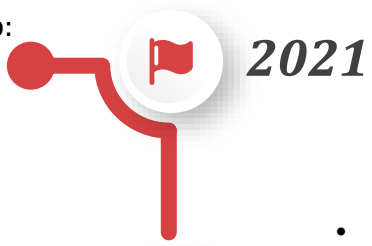


INTRODUCTION

Goal

Invest US\$50 million to:

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



2019 - 2020



- 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



First expansion
(175 new stations)

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



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:

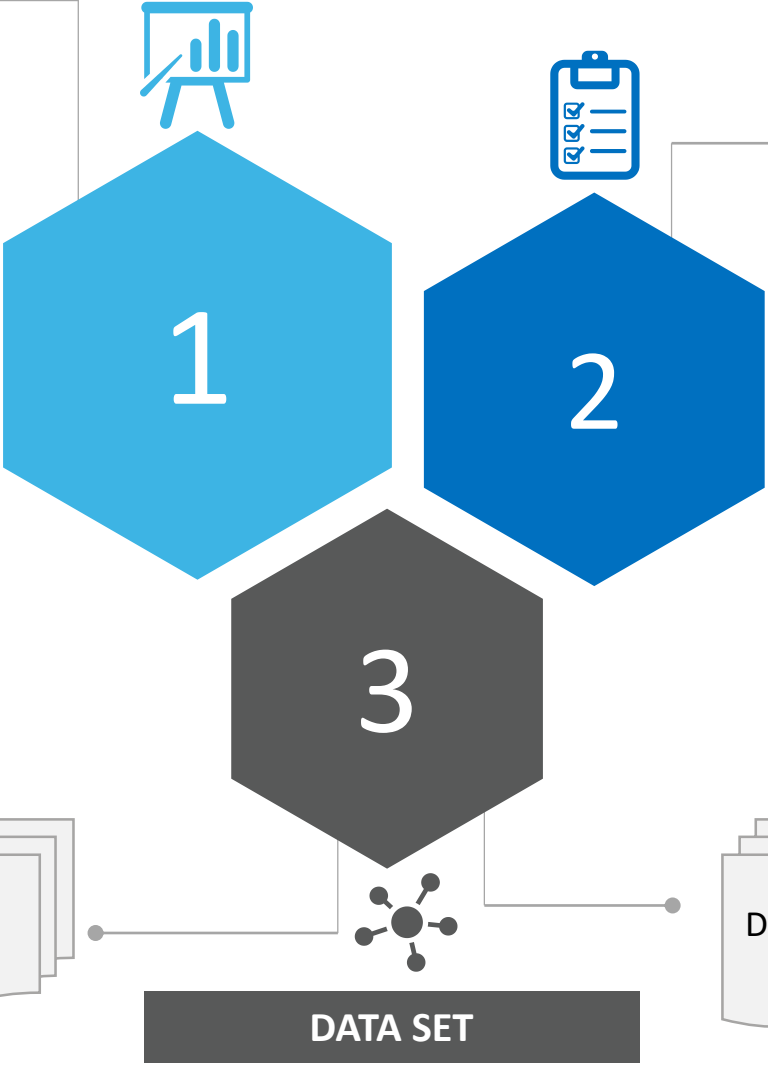
- Provide methodologies and various tools used in the process
- Provide data analysis and visualization
- Put forward a future state blueprint for the new stations and bikes allocation process



PROPOSED FINDING

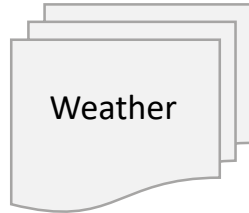
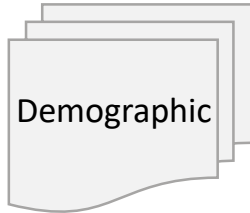
Our final deliverables will enable Divvy leadership to:

- 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 SET



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



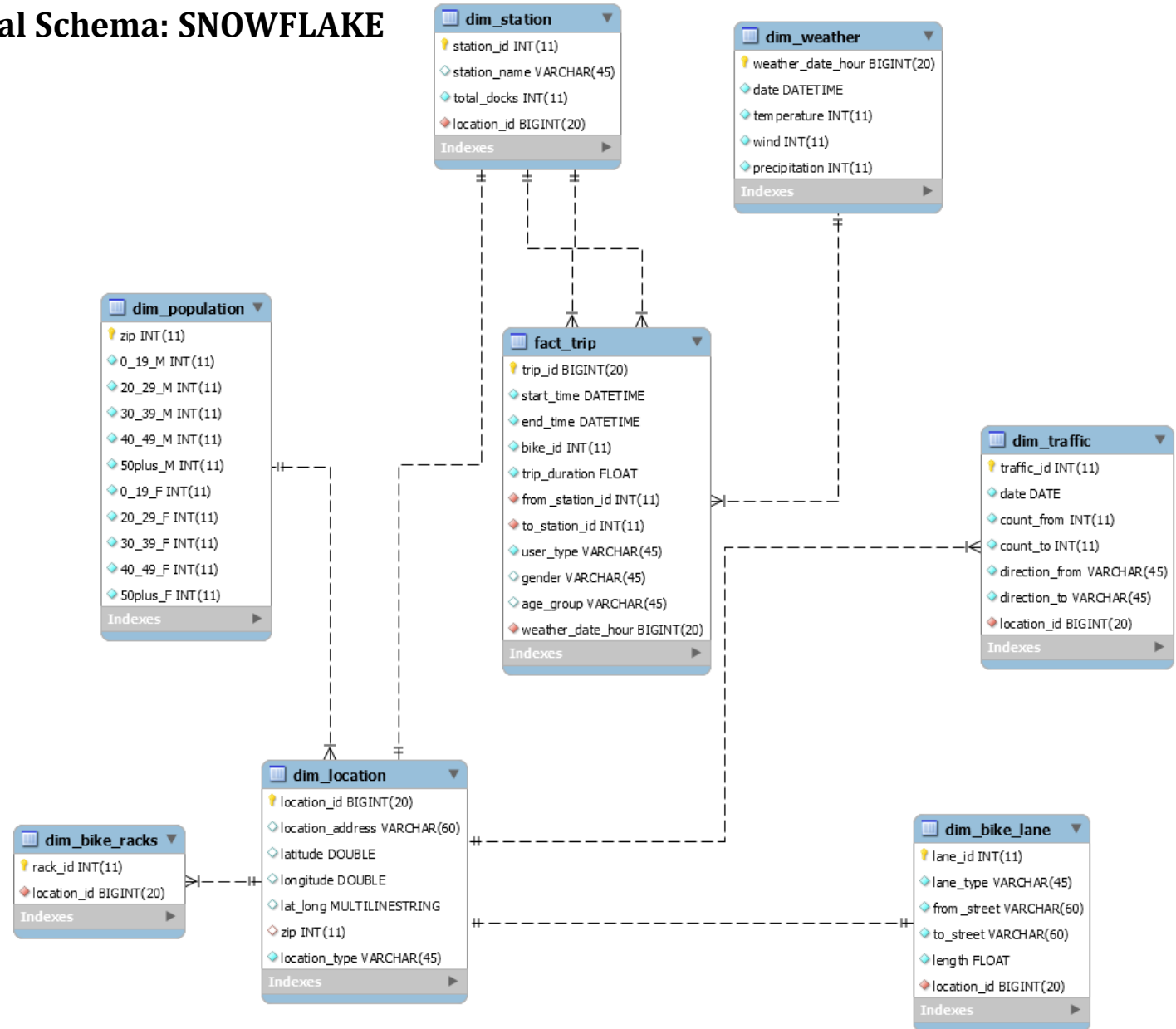
Dimensional Schema: SNOWFLAKE

DDL

```

-- MySQL Workbench Forward Engineering
1
2
3 SET SQL_MODE=NO_AUTO_VALUE_ON_ZERO;
4 SET SQL_MODE=NO_AUTO_VALUE_ON_ZERO;
5
6
7
8
9
10
11
12
13
14
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
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
    
```

DML





1

Data

DIY



Data from various sources

DIY

2

ETL Process

Refine



3

Data Warehouse



Google Cloud Platform

4

Data Analysis



5

Data Visualization



6

Findings and Recommendations



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,
COUNT(trip_id) AS NoOfTrips
FROM fact_trip
GROUP BY TimeOfDay, DateType
ORDER BY TimeOfDay DESC;
    
```

```

# Number of TripIn per zip
• SELECT
  dl.zip,
  COUNT(ft.to_station_id) AS TripIn
FROM
  fact_trip ft
INNER JOIN dim_station ds ON ds.station_id = ft.to_station_id
LEFT JOIN dim_location dl ON dl.location_id = ds.location_id
GROUP BY zip
ORDER BY TripIn DESC;
    
```



- 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...



- 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	Gender	Age Group
1	Rogers Park	Male	0-19
2	West Ridge	Female	20-29
3	Uptown		
4	Lincoln Square		
5	North Center		
6	Lake View		
7	Lincoln Park		
8	Near North Side		
9	Edison Park		
10	Norwood Park		
11	Jefferson Park		
12	Forest Glen		
13	North Park		
14	Albany Park		
15	Portage Park		
16	Irving Park		
17	Dunning		
18	Montclare		
19	Belmont Cragin		
20	Hermosa		
21	Ancrudae		
22	Logan Square		
23	Humboldt Park		
24	West Town		
25	Austin		



- 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.

```

1 install.packages("rpopr")
2 library(rpopr)
3
4 datapath = "D:/Users/Somey/Tammy/ChicagoData/Engineering_Platform/ProdProject/Processing/"
5
6 #Rack1
7 rack1 = read.csv(filepathdatapath, "rack1.csv", sep=";")
8 #Rack2
9 rack2 = read.csv(filepathdatapath, "rack2.csv", sep=";")
10 #Rack3
11 rack3 = read.csv(filepathdatapath, "rack3.csv", sep=";")
12 write.table(rack1, file = paste(datapath, "rack1.csv", sep = "/"), row.names = F)
13
14 # Average distance per station
15 #Rack1
16 #Rack2
17 #Rack3
18 # Average distance per station
19 #Rack1
20 #Rack2
21 #Rack3
    
```



Data Analysis and Visualization





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 FROM
13     (SELECT
14         ds.station_id,
15         ds.station_name AS stationName,
16         ds.total_docks AS totalDocks,
17         HOUR(ft.start_time) AS TimeOfDay,
18         COUNT(ft.from_station_id) as tripFrom
19     FROM
20         fact_trip ft
21         INNER JOIN
22         dim_station ds ON ds.station_id = ft.from_station_id
23     GROUP BY
24         ds.station_id, TimeOfDay
25     ORDER BY
26         ds.station_id,TimeOfDay ASC) AS TripFrom
27     INNER JOIN
28     (SELECT
29         ds.station_id,
30         ds.station_name AS stationName,
31         ds.total_docks AS totalDocks,
32         HOUR(ft.end_time) AS TimeOfDay,
33         COUNT(ft.to_station_id) as tripTo
34     FROM
35         fact_trip ft
36         INNER JOIN
37         dim_station ds ON ds.station_id = ft.to_station_id
38     GROUP BY
39         ds.station_id, TimeOfDay
40     ORDER BY ds.station_id,TimeOfDay ASC) AS TripTo ON TripFrom.station_id = TripTo.station_id
41 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 FROM
102     (SELECT
103         ds.station_id,
104         ft.trip_id,
105         dl.latitude,
106         dl.longitude
107     FROM
108         dim_location dl
109         INNER JOIN
110         dim_station ds ON dl.location_id=ds.location_id
111         INNER JOIN
112         fact_trip ft ON ds.station_id=ft.from_station_id) AS FrS
113     INNER JOIN
114     (SELECT
115         ds.station_id,
116         ft.trip_id,
117         dl.latitude,
118         dl.longitude
119     FROM
120         dim_location dl
121         INNER JOIN
122         dim_station ds ON dl.location_id=ds.location_id
123         INNER JOIN
124         fact_trip ft ON ds.station_id=ft.to_station_id) AS TrS ON FrS.trip_id=TrS.trip_id
125 WHERE
126     FrS.station_id != TrS.station_id;
```

Customer Profiling

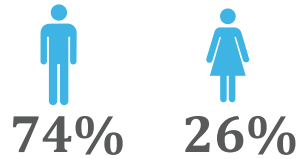


Users Type

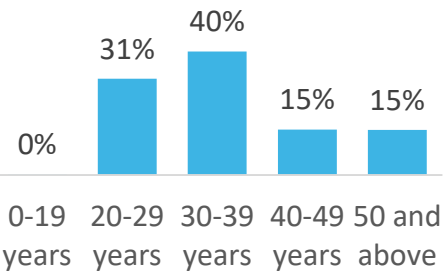


- Subscriber
- Non-subscriber

Gender



Age Group



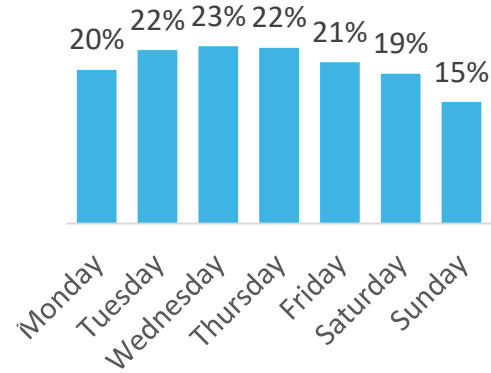
Average Distance Travel



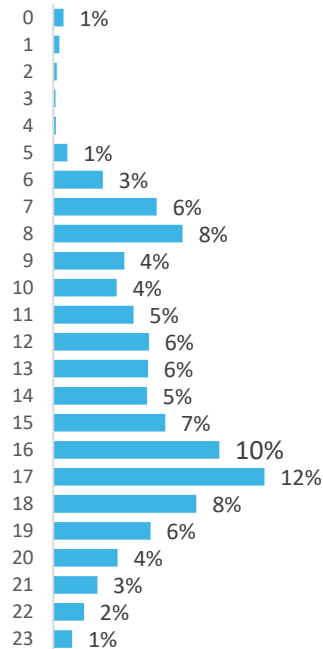
Average Trip Duration



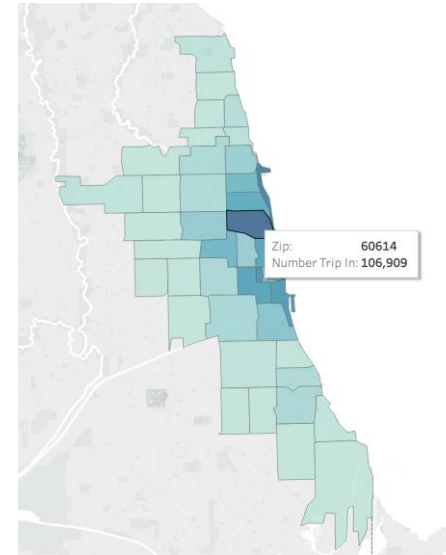
Trip by day



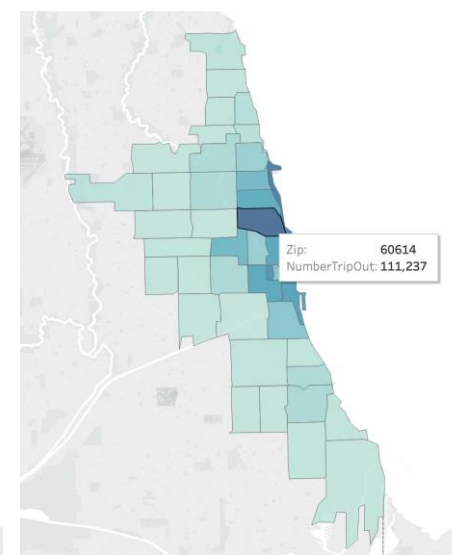
Trip by hour



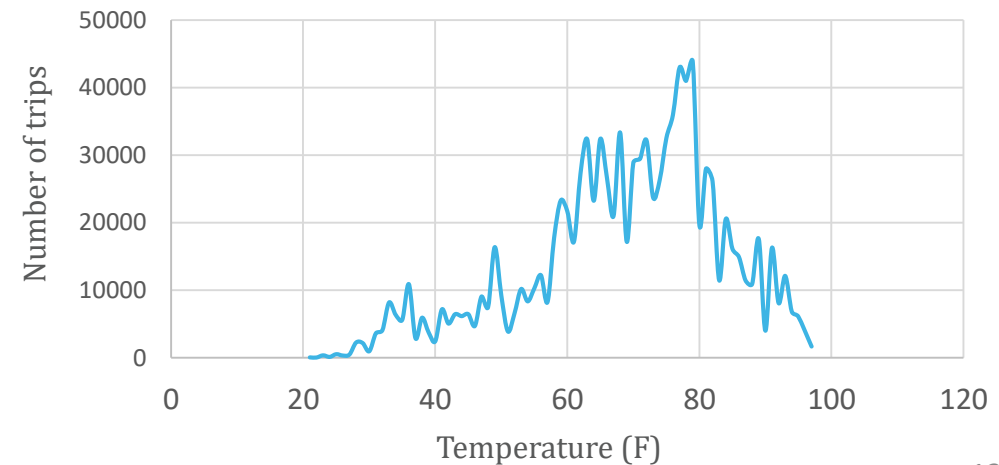
Trip-in by area



Trip-out by area



Trip by weather



Findings by zip code



Zipcode Analysis	Demographic													Traffic Vehicle Volume	BIKE Number of Bike racks	Number of stations	Dewy Stations Total # of docks	Avg # of docks	Avg of Avg Distance of stations from other stations (miles)	Dewy Trips				Subscriber %
	Population Total	Gender		Age																Trips Out	Trips In	Net	% of total trips	
		Male%	Female%	0_9	10_19	20_29	30_39	40_49	50plus	0_9	10_19	20_29	30_39											
60605	24672	48.5%	51.5%	0.0678	0.138	0.1094	0.0585	0.1719	0.1451	0.1783	0.0603	0.1206	8,300	358	13	550	28.9	4.75	68,902	65,243	(3,659)	6.33%	60.72%	
60601	11195	49.42%	50.58%	0.0639	0.1702	0.0933	0.0595	0.1703	0.1746	0.0933	0.0571	0.173	23,800	1	11	364	11.2	4.53	68,506	63,594	(4,912)	5.22%	72.45%	

Sample zip code findings

Zip code	60605
Total Population	24672
Male%	0.4852
Female%	0.5148
0_19_M	0.0678
20_29_M	0.138
30_39_M	0.1094
40_49_M	0.0565
50plus_M	0.1134
0_19_F	0.0719
20_29_F	0.1451
30_39_F	0.1169
40_49_F	0.0603
50plus_F	0.1206
Vehicle Volume	8,300
Number of Bike racks	356
Number of stations	19
Total # of docks	550
Avg # of docks	28.94736842
Avg of Avg Distance of stations from other stations (miles)	4.754563244
Trips Out	68,902
Trips In	65,243
Net	-3,659
% of total trips	0.063305521
Subscriber %	0.606542626

Sample calculation

* Sample calculation
- Distance between stations

Zipcode	Male%	Female%	0_9	10_19	20_29	30_39	40_49	50plus	Vehicle Volume	Number of Bike racks	Number of stations	Total # of docks	Avg # of docks	Avg of Avg Distance of stations from other stations (miles)	Trips Out	Trips In	Net	% of total trips	Subscriber %				
60605	0.4852	0.5148	0.0678	0.138	0.1094	0.0565	0.1134	0.0719	0.1451	0.1783	0.0603	0.1206	8300	356	13	550	28.9	4.75	68902	65243	-3659	6.33%	60.72%

- Latitude & Longitude for 608 existing stations
- Used complex formula involving trigonometry to find distances
 - $=IFERROR(6371*ACOS(COS(RADIANS(90-B63)))*COS(RADIANS(90-VS$2))+SIN(RADIANS(90-$B$63))*SIN(RADIANS(90-VS$2)))/1.609,0)$
- Resulting in over 30k distance values for each pair of stations

Factors considered for analysis. For each zip code we found:

- Total population
- Male & Female %
- %s of different age_groups
- Vehicle volume
- Number of bike racks
- Number of stations
- Number of docks
- Avg of avg distance of stations from other stations
- Trip Out, Trip In, Net
- % of subscribers





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.

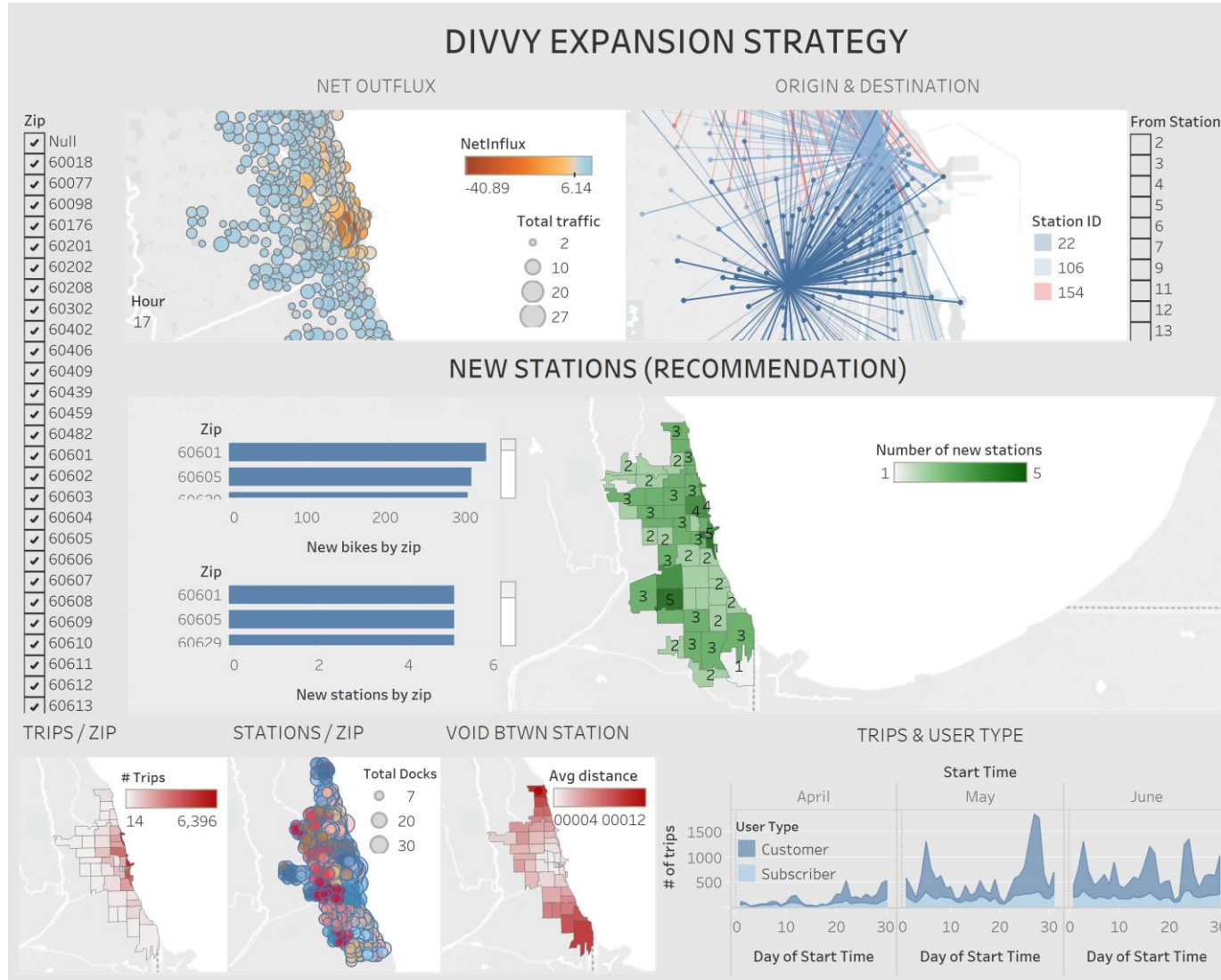
Our scoring methodology

- 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%	-



Tableau Visualization



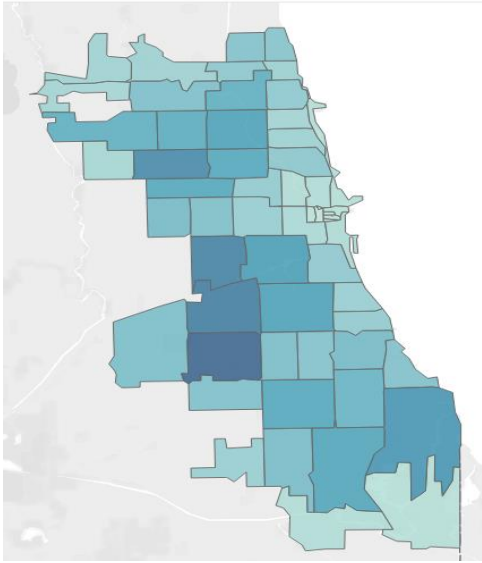
Derived recommendation from trip and zip demographics:

- **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 Btwm 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')

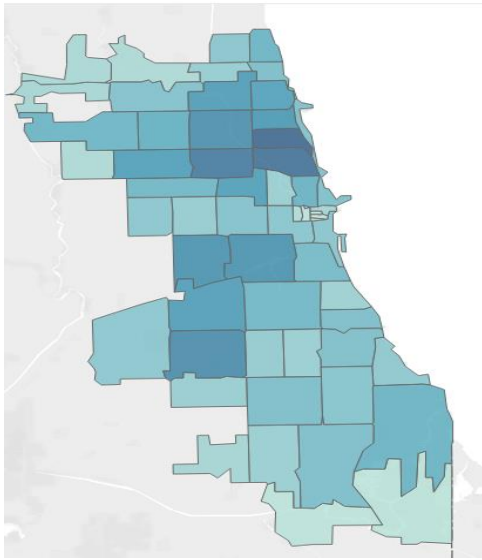
Demographics by Zip Code



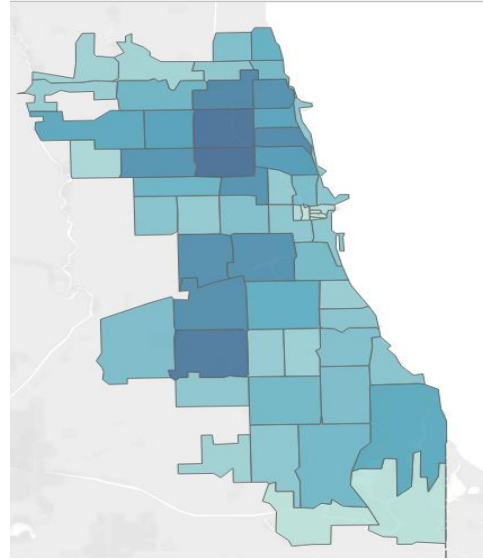
0-19 ZIP



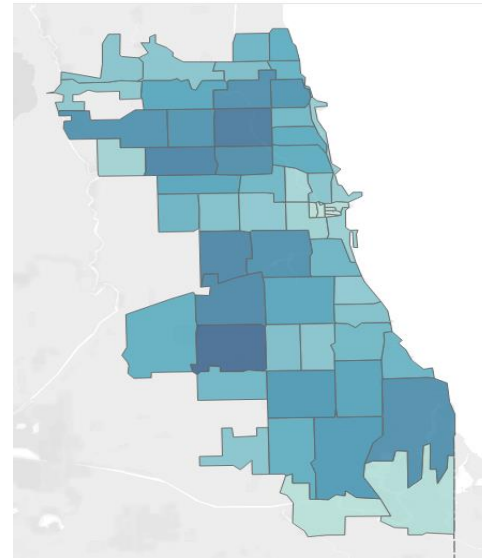
20-29



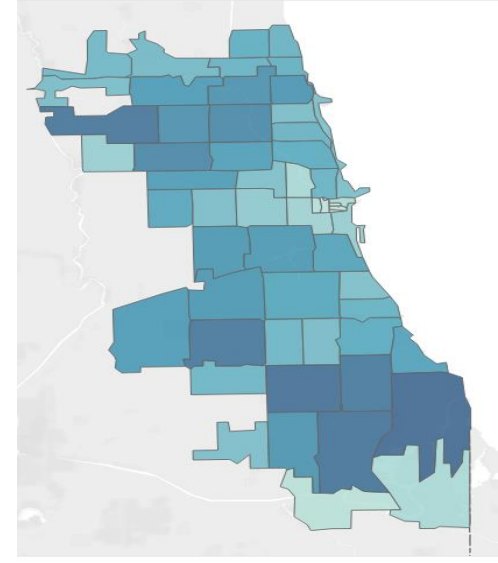
30-39



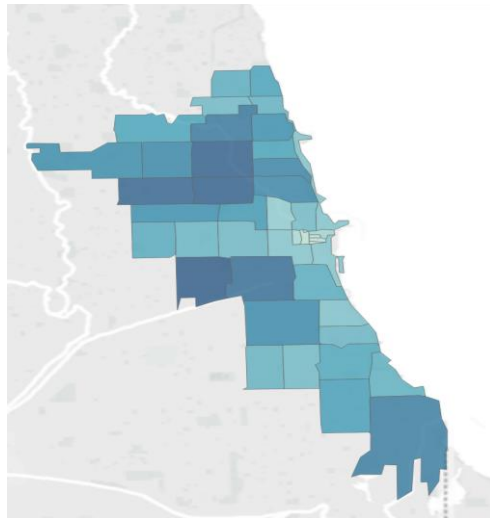
40-49



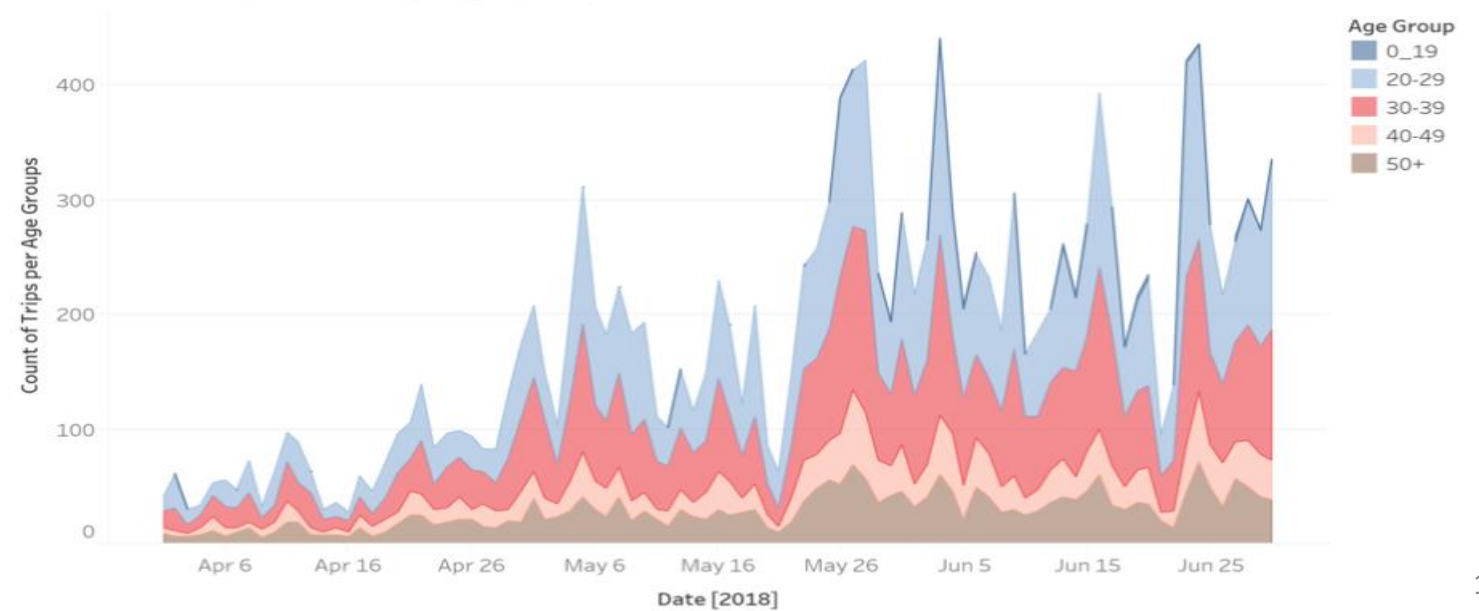
50+



Total population by zip



Number of trips taken by age groups





Summary



Recommendations and Future Vision:

- Increase stations in ZIPs further 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!
