

A Data-driven Approach to Eliminating Pedestrian Collisions in Toronto

Jason Kim

jason2.kim@ryerson.ca

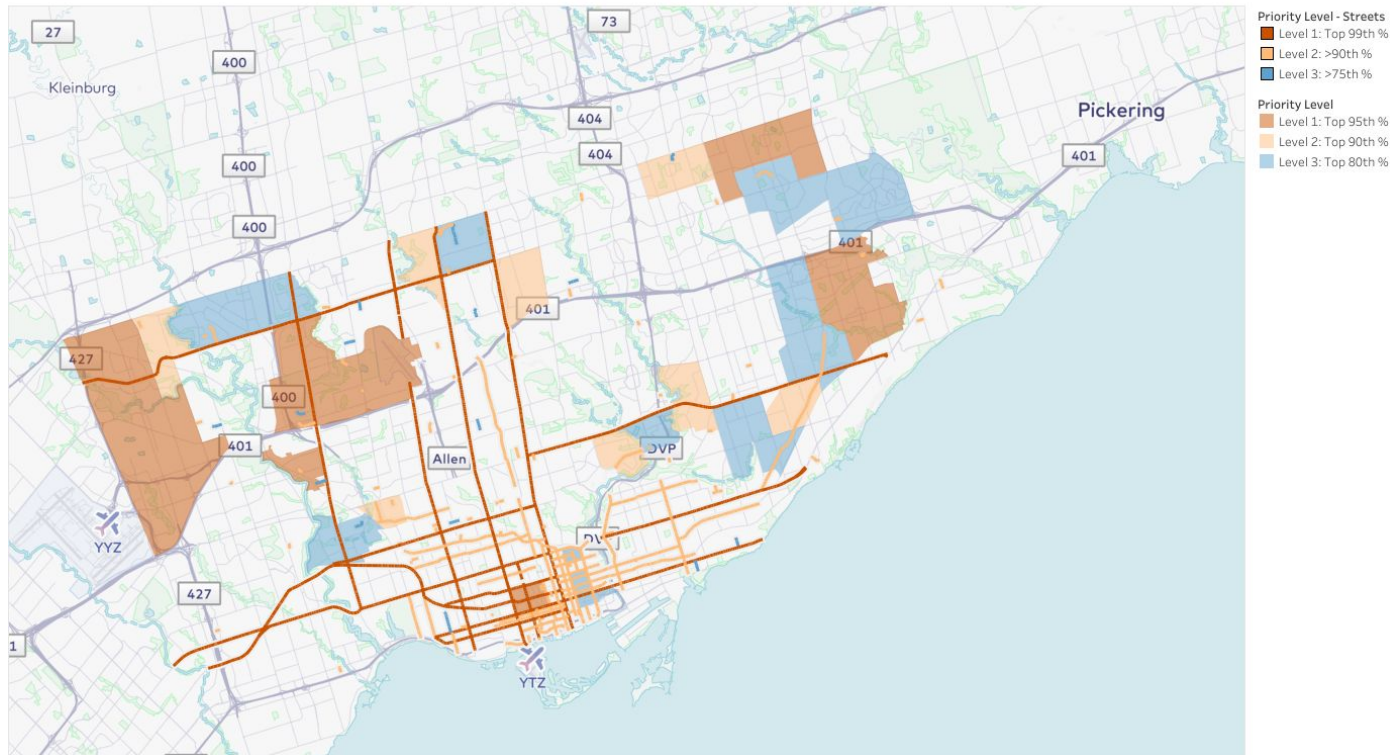
<https://github.com/jasonukim/capstone-repo>

Context

- Just a few months ago, the City of Toronto approved \$22 million in additional funding for accelerating its Vision Zero Road Safety Plan
- The goal of this plan is to eliminate pedestrian deaths from car collisions by 2021
- 2018 is one of the worst years for pedestrian and cyclist fatalities, exceeding 2017's peak year of 162 deaths or major hospitalizations
- Politically sensitive topic due to perceptions:
 - “War against the car”
 - Collisions are caused by irresponsible pedestrians rather than drivers

Previous Work

Toronto High Injury Network



Research Question

- To identify and prioritize zones within Toronto that have a high risk of pedestrian collisions so that safety improvements could be more efficiently implemented by the City

Methodology

- Joined multiple datasets to the collision dataset to increase level of detail about factors that affect collision density
- Used various clustering algorithms in order to identify high collision zones
 - Kernel Density Estimation (KDE)
 - K-means Clustering
 - Density-based Spatial Clustering of Applications with Noise (DBSCAN)
- Used Average Silhouette to test the “goodness of fit” for clusters
- Random Forest to classify unseen collisions into clusters and check Variable Importance

Tools Used

- R
 - General: Caret, corrplot, sqldf, mlbench, MASS, sp
 - Mapping: Maptools, ggmap, rgeos
 - Visualization: Ggplot2, knitr
 - Clustering, Silhouette, and Classification: Factoextra, NbClust, dbscan, cluster, RandomForest
 - Parallel Processing: doSNOW
- QGIS
 - For spatial joins on shapefiles (a special kind of dataset containing geo data)

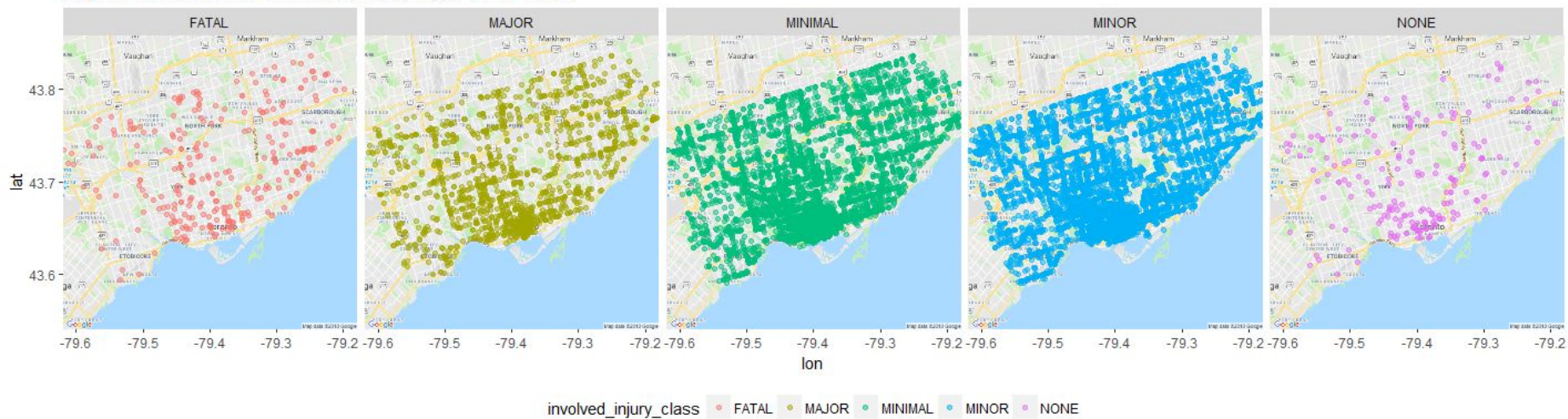
Dataset Details

- Over 10,000 rows of collisions
- 88 attributes from 9 different open data sources
- Geospatial, demographic, economic, health, language, SES, etc.
- Full data dictionary:

<https://github.com/jasonukim/capstone-repo/blob/master/Datasets/Data%20Dictionary.md>

Exploratory Data Analysis

Reported Collisions in Toronto by Injury Type (2007-2017)



Neighbourhood	Collisions	Street	Collisions
Waterfront Communities-The Island (77)	608	YONGE ST	555
Bay Street Corridor (76)	567	DUNDAS ST W	371
Church-Yonge Corridor (75)	367	BATHURST ST	366
Downsview-Roding-CFB (26)	286	BLOOR ST W	326
Islington-City Centre West (14)	286	EGLINTON AVE E	319

Road Type Collisions

MAJOR ARTERIAL	10318
MINOR ARTERIAL	2645
COLLECTOR	1294
LOCAL	1276

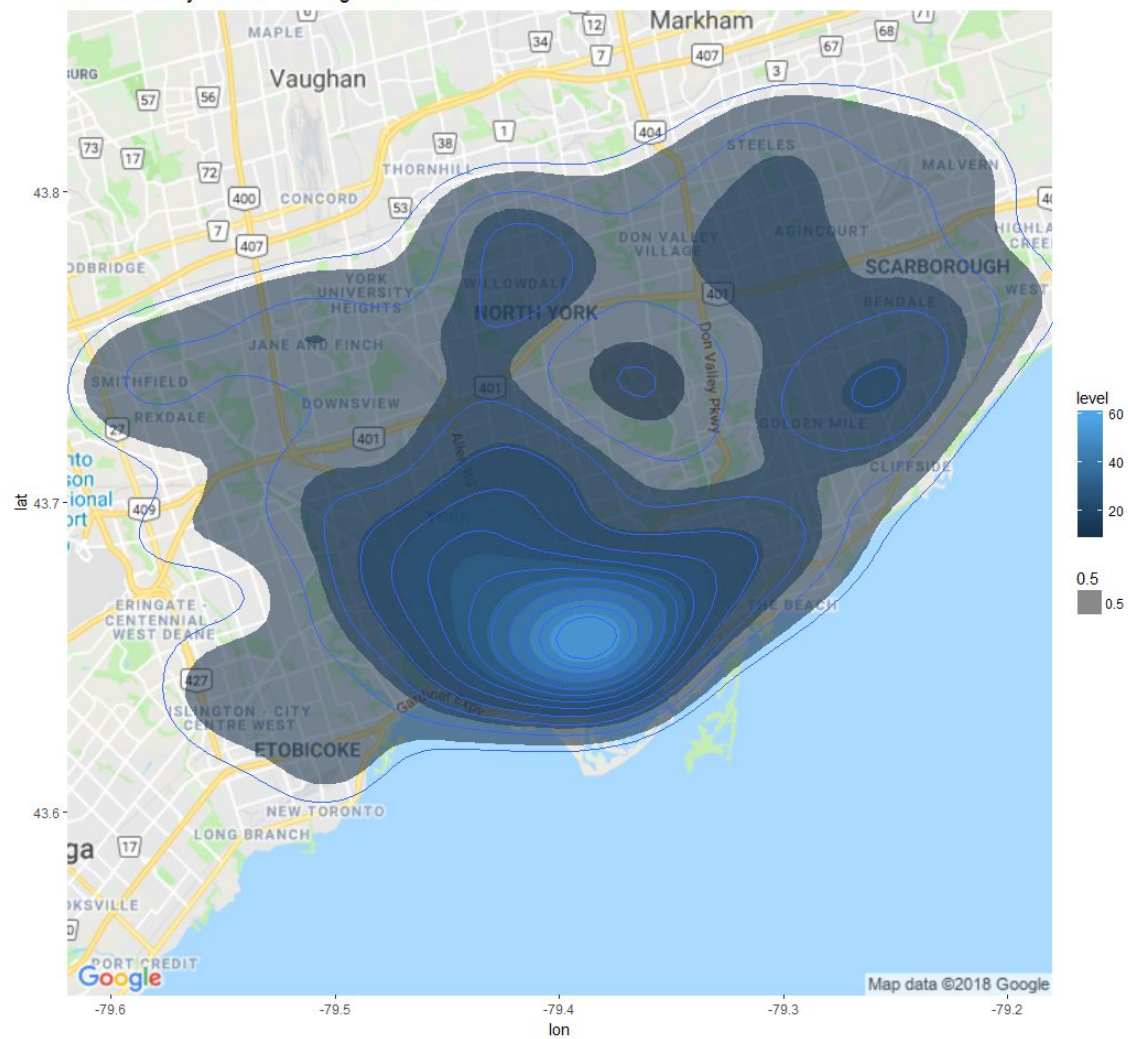
Location Collisions

INTERSECTION	12186
MID-BLOCK	3304

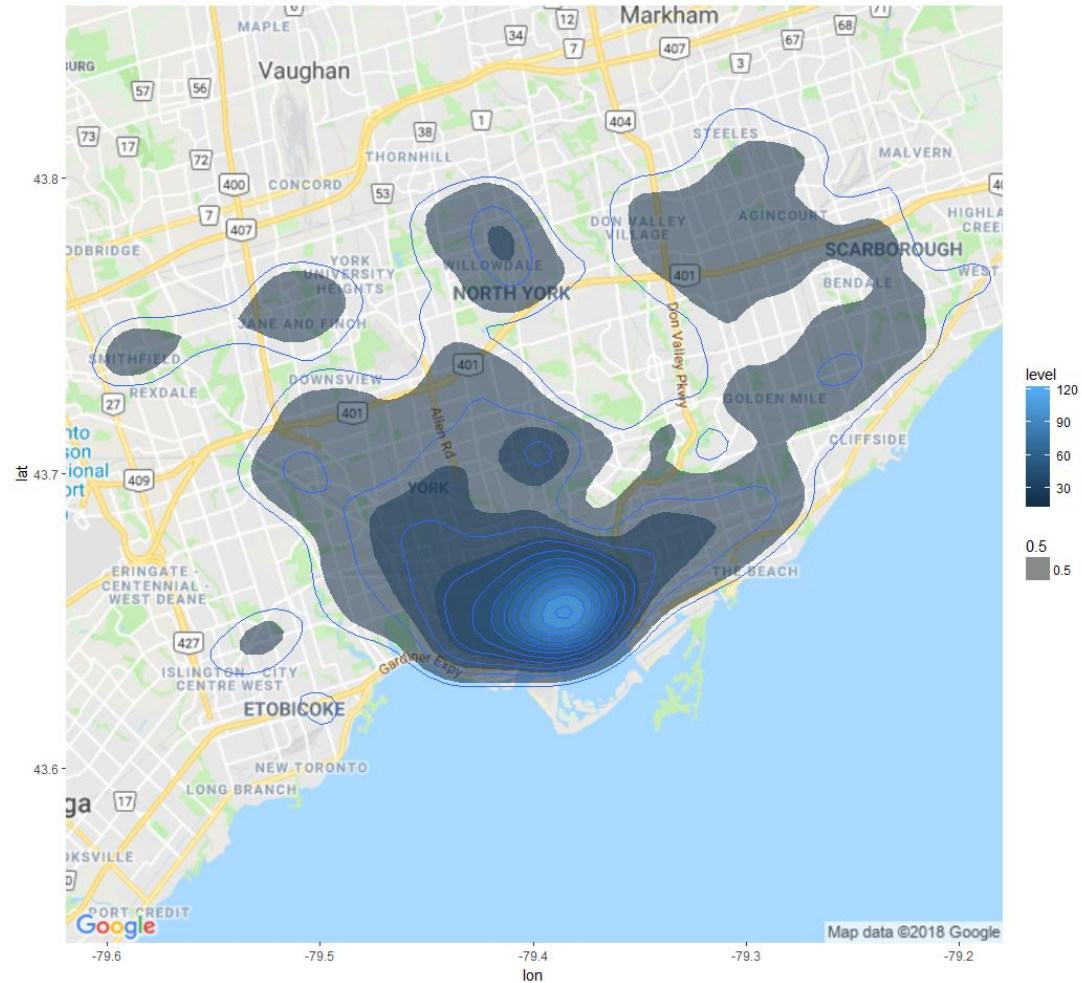
Visibility Condition Collisions

CLEAR	12311
RAIN	2563
SNOW	457
OTHER	72
FREEZING RAIN	41
FOG, MIST, SMOKE, DUST	40
DRIFTING SNOW	29

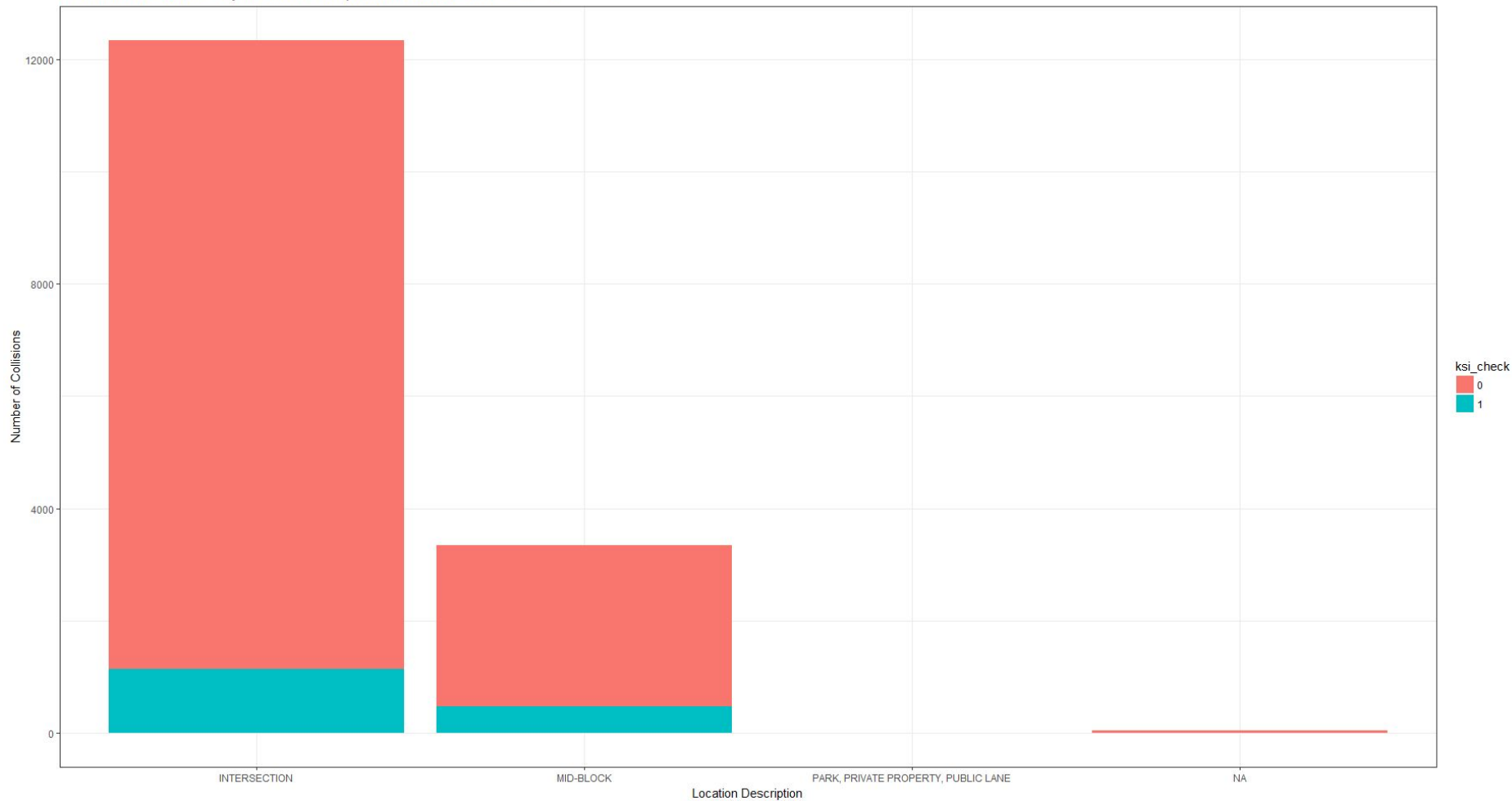
Kernel Density Estimation of High KSI Zones



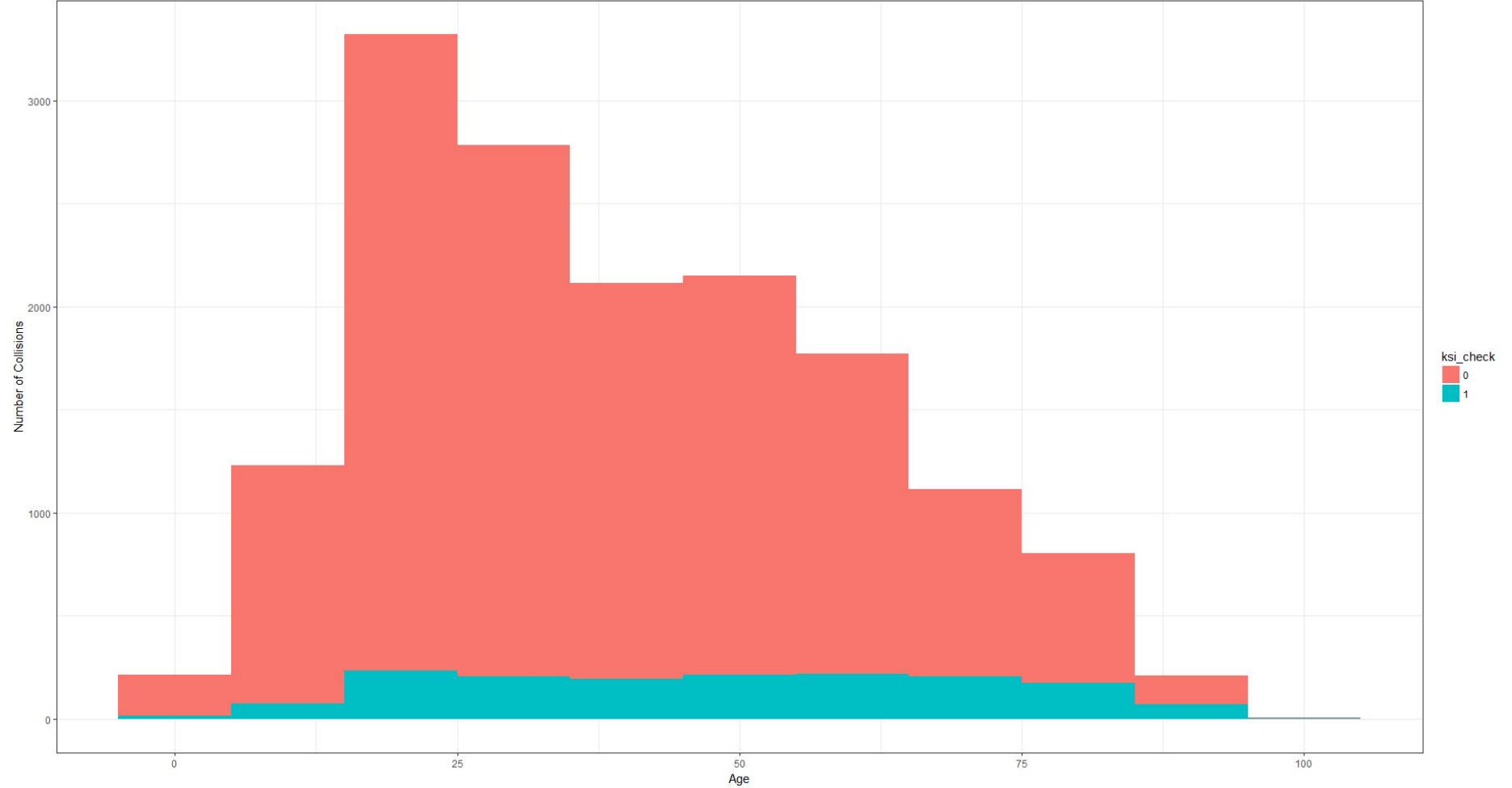
Kernel Density Estimation of High Non-KSI Zones



Distribution of Collisions by Location Description Across KSIs and non-KSIs

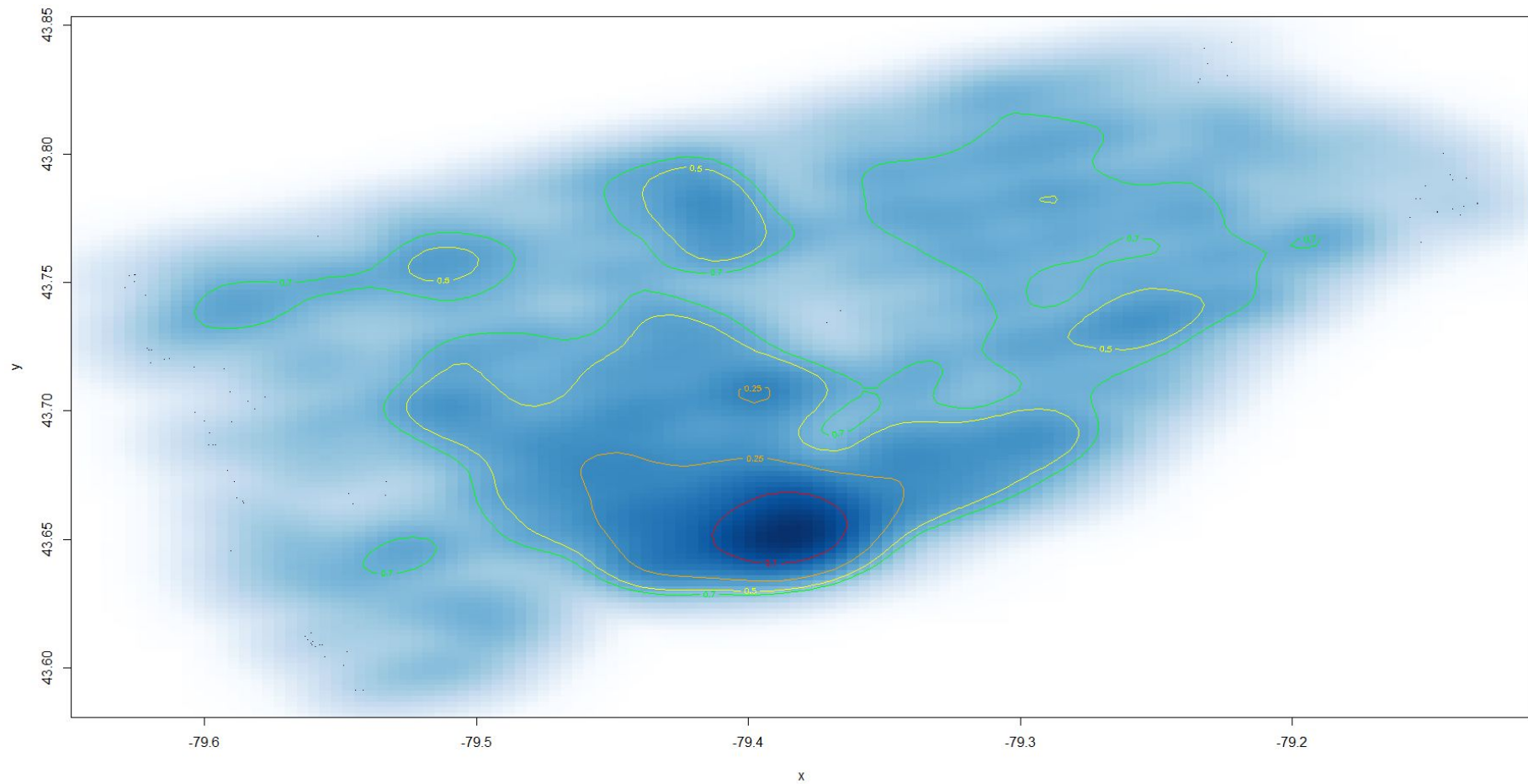


Distribution of Collisions by Age Across KSI vs. Non-KSI Collisions

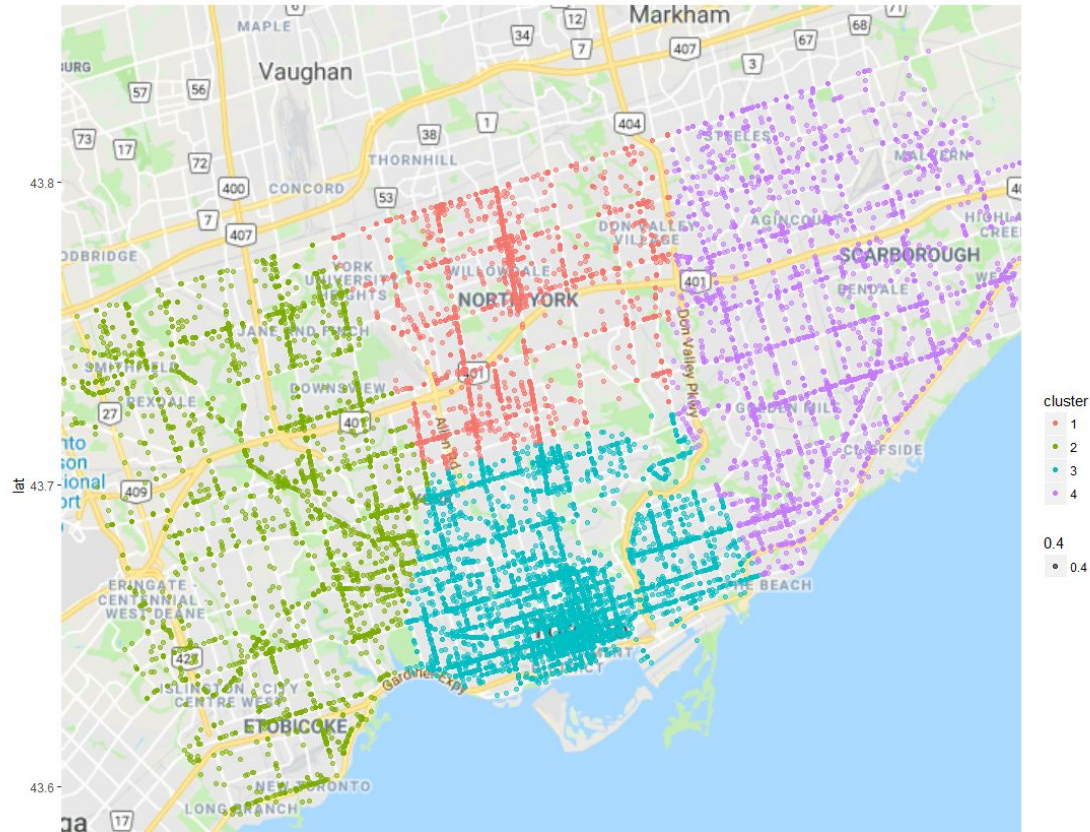


Clustering by Collision Location Only

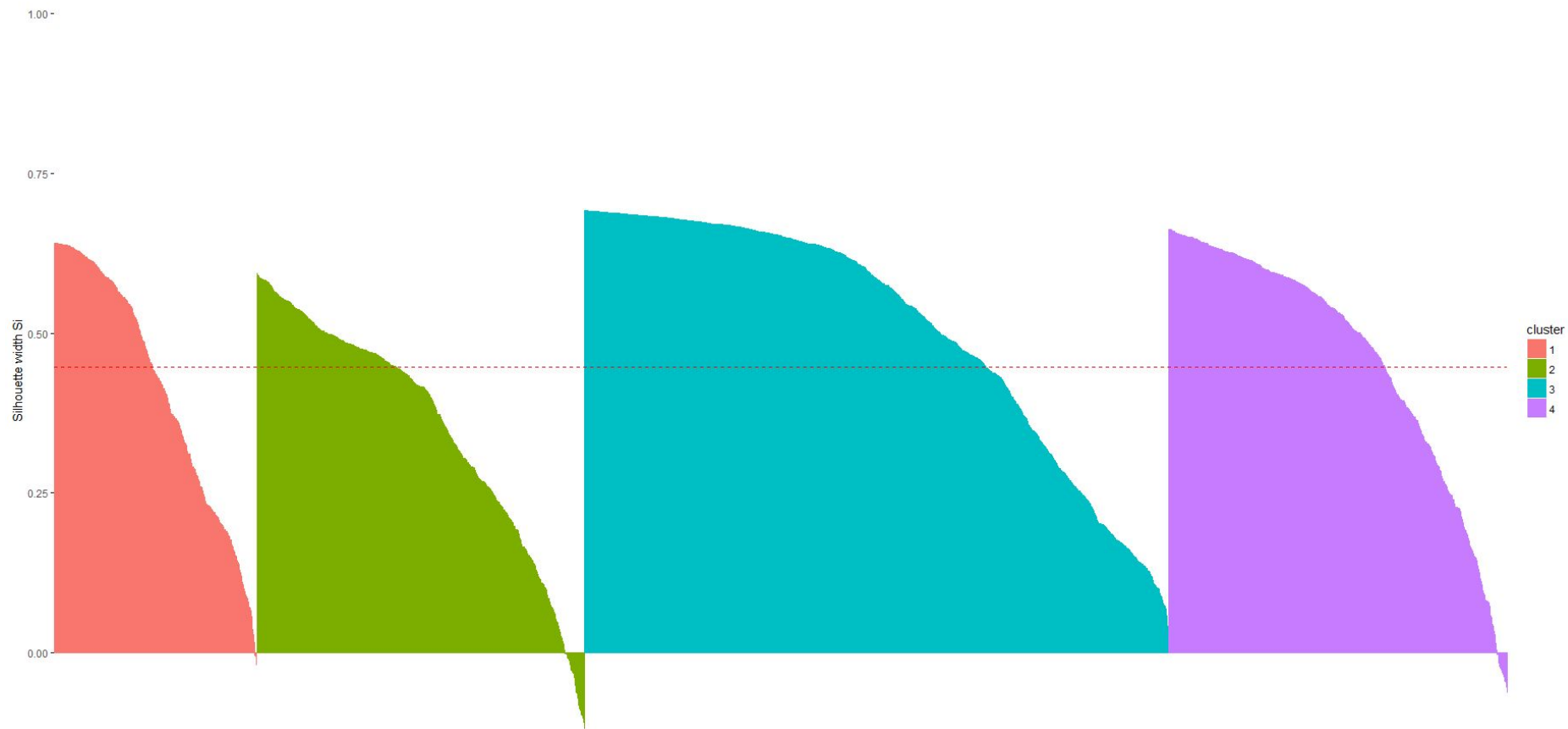
KDE



K-Means Clustering

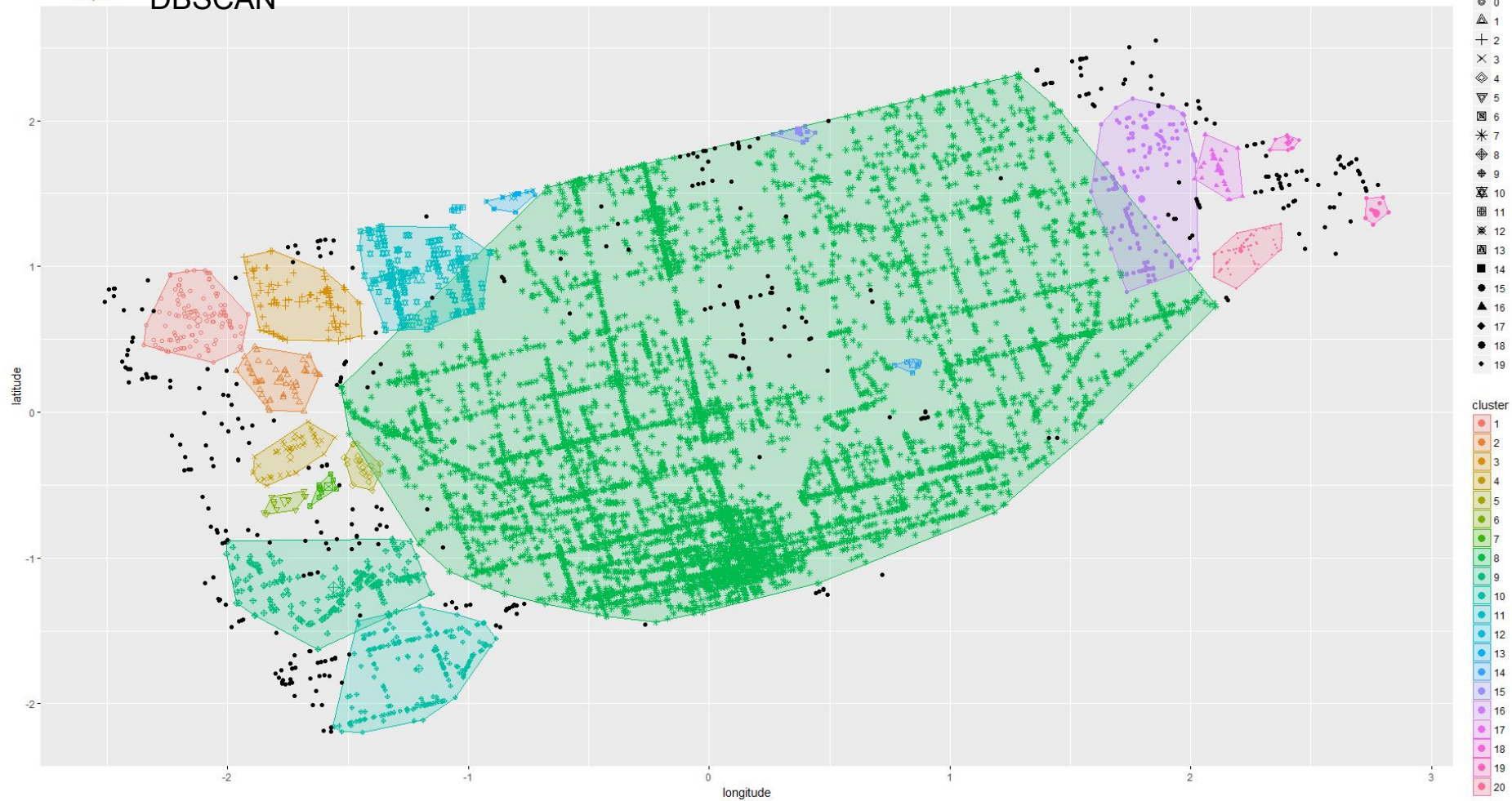


Clusters silhouette plot
Average silhouette width: 0.45

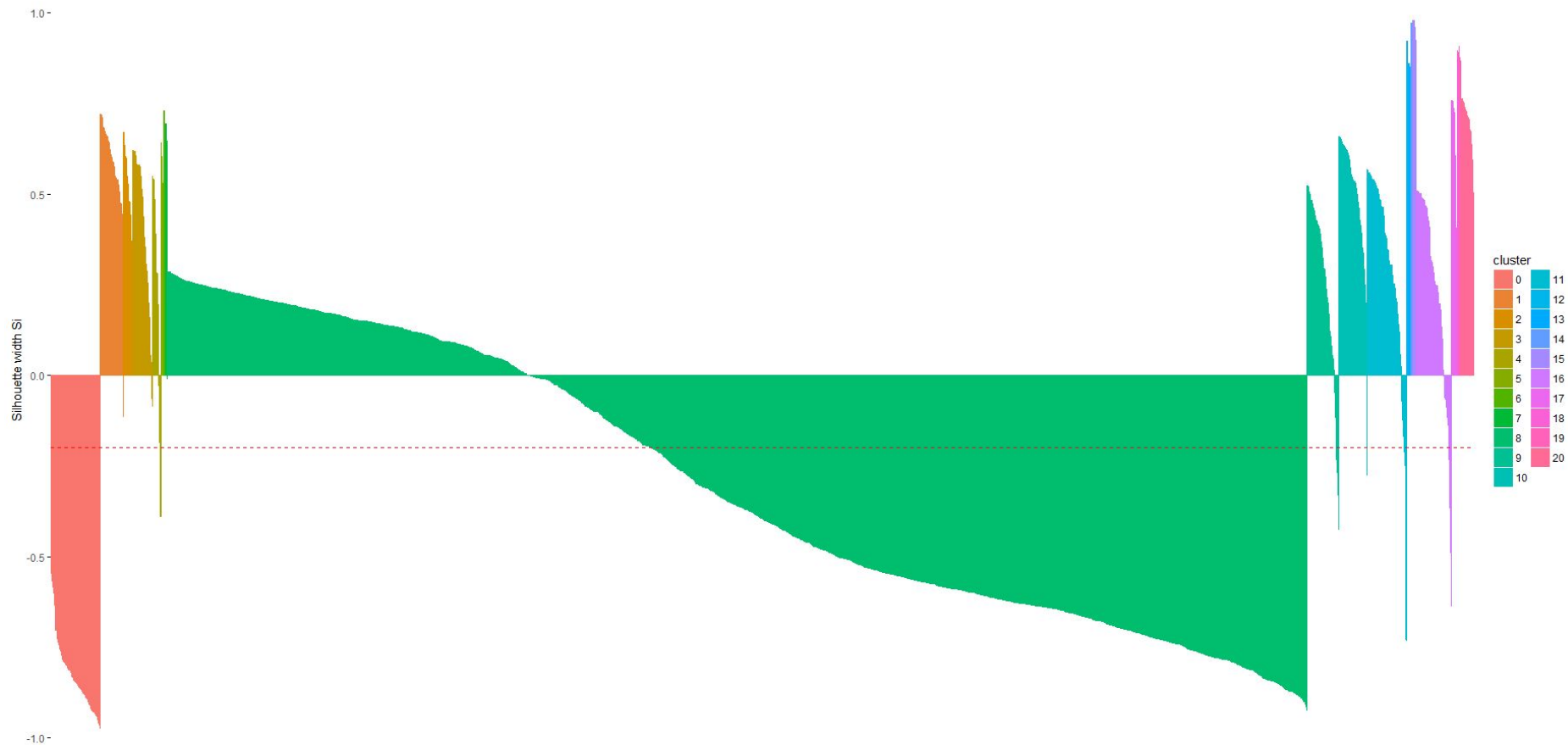


Cluster plot

DBSCAN

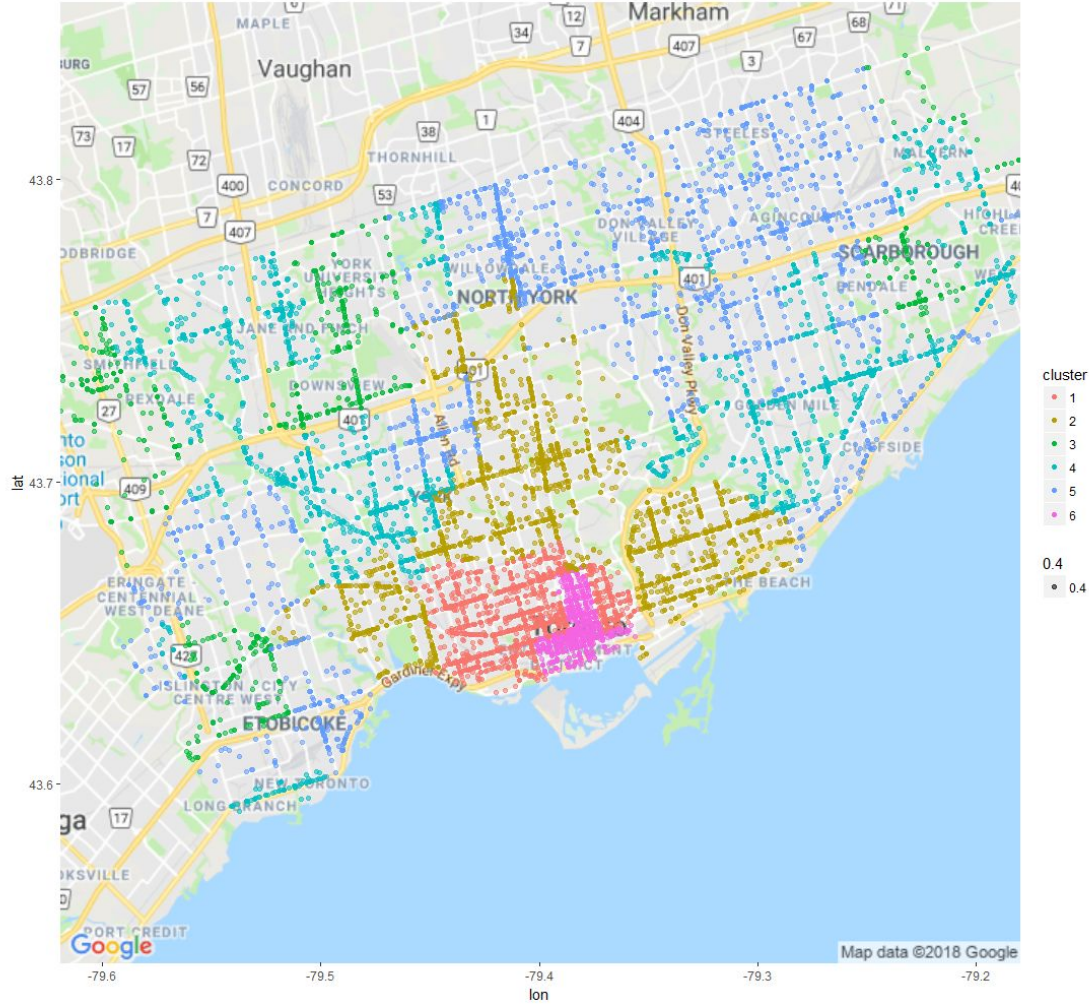


Clusters silhouette plot
Average silhouette width: -0.2

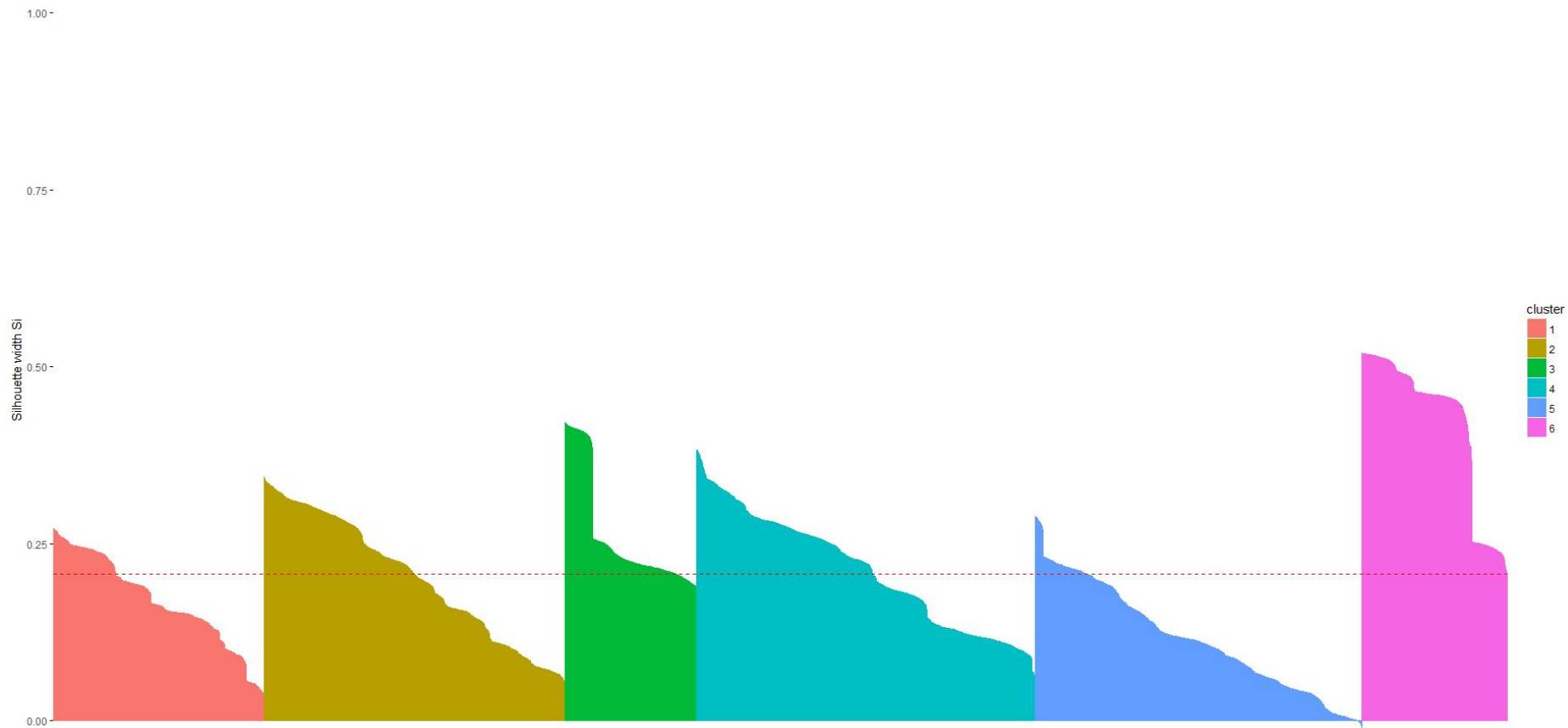


Clustering by Normalized Numerical Variables

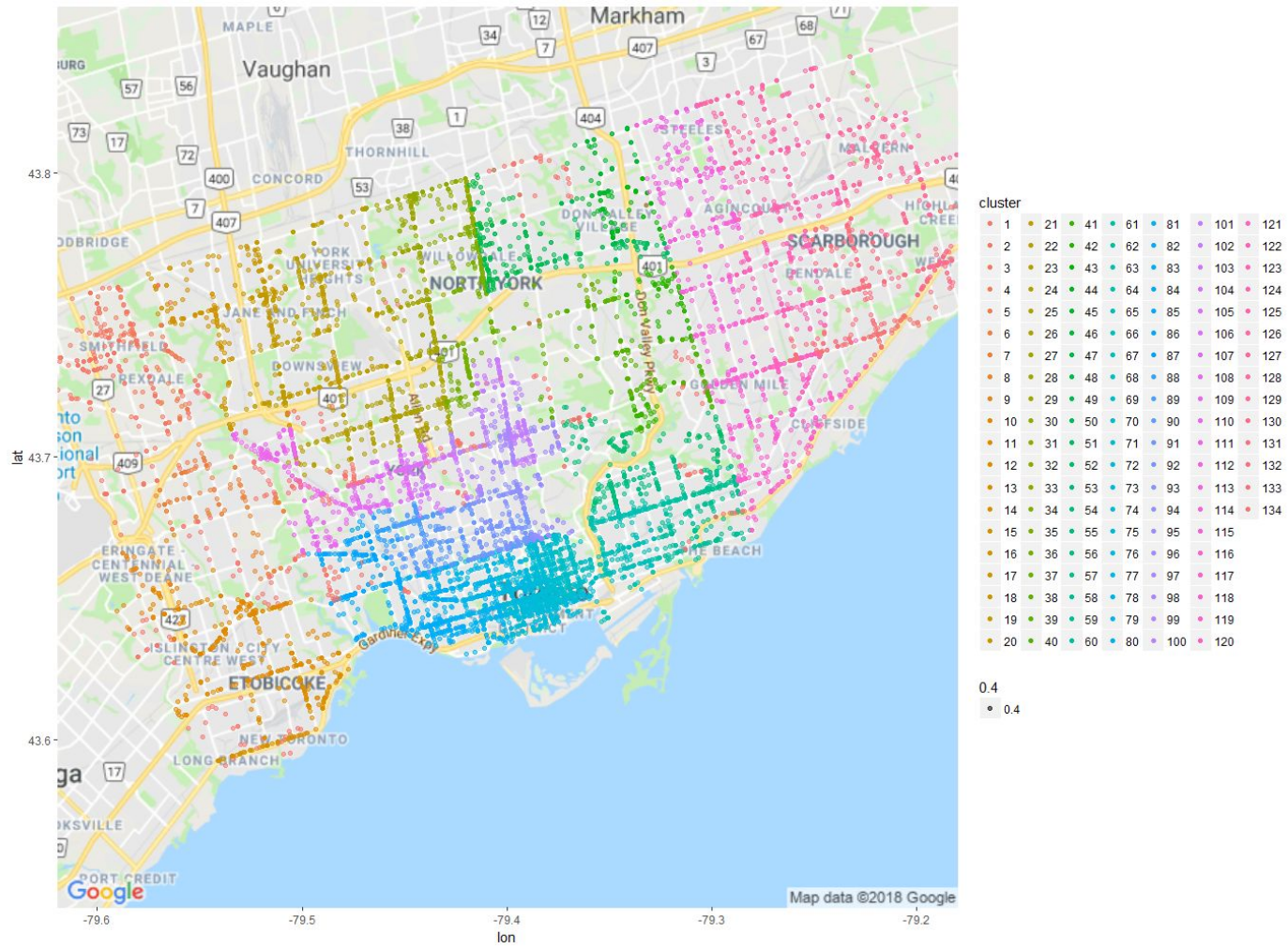
K-means



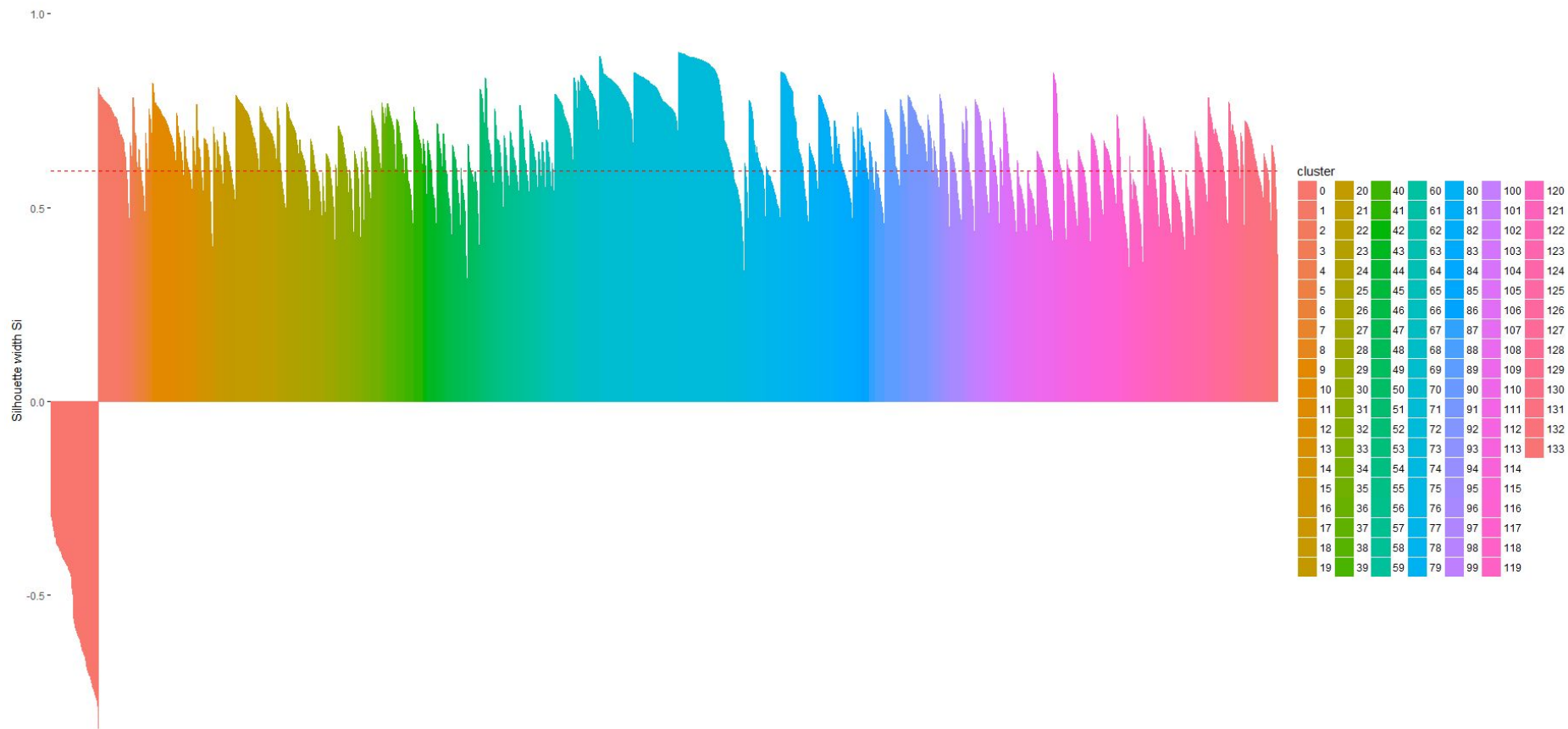
Clusters silhouette plot
Average silhouette width: 0.21



DBSCAN

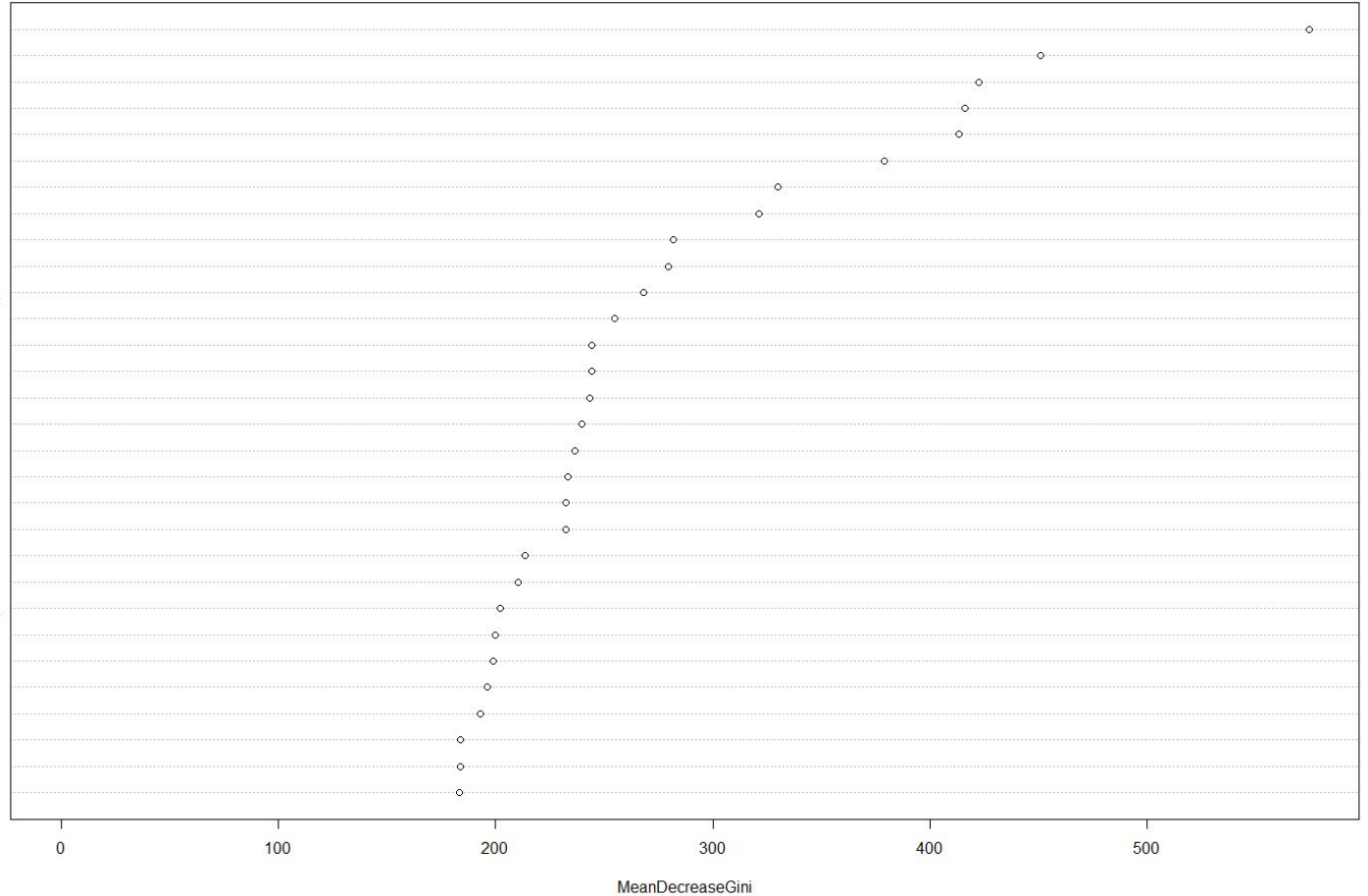


Clusters silhouette plot
Average silhouette width: 0.59



rf_model

Number.of.Businesses
 Population.2016
 City.Grants.Funding..
 Local.Employment
 Pedestrian.Other.Collisions
 TTC.Stops
 Road.Kilometres
 Working.Age..25.54.years.
 Traffic.Collisions
 Land.area.in.square.kilometres
 X..Time.leaving.for.work...Between.8.a.m.and.8.59.a.m.
 Watermain.Breaks
 Diversity.Index...2008
 X..Commute.in.Walked
 Voter.Turnout...2008
 Social.Assistance.Recipients
 Salvation.Army.Donors
 Debt.Risk.Score
 Walk.Score
 Child.Care.Spaces
 Commuting.duration...15.to.29.minutes
 X..Visible.Minority
 X..Time.leaving.for.work...Between.6.a.m.and.6.59.a.m.
 Children...0.14.years.
 Population.density.per.square.kilometre
 Unemployment.rate
 Commuting.duration...45.to.59.minutes
 Seniors...65..years.
 Road.Volume
 Neighbourhood.Equity.Score



Conclusion

- KDE is very useful for clustering densities, but not much else; KDE is the dominant method used in traffic research
- DBSCAN is incredibly strong at clustering collisions with shared characteristics that don't follow a circular pattern
- Although RF identified which variables are important as a whole, not sufficient to tell us about each unique cluster
- Unsupervised learning makes supervised learning even more powerful and vice versa
- The City should adopt a clusters-based approach rather than street or neighbourhood-level due to their size

A Data-driven Approach to Eliminating Pedestrian Collisions in Toronto

Jason Kim

jason2.kim@ryerson.ca

<https://github.com/jasonukim/capstone-repo>