

Applied Data Science Capstone - The Battle of the Neighbourhoods - Final Report

Understanding COVID-19 in Toronto Neighbourhoods – A Data-driven Approach For Vaccination Prioritization

By Chris Blackwood

Introduction

COVID-19 has been a major health and economic issue worldwide since early 2020. The virus was first identified in humans in Wuhan, China in December of 2019, and the World Health Organization recognized a global pandemic in March of 2020 [1]. As of late January 2021, over 100 million cases have been identified worldwide, and over 2 million people have died [2].

Many organizations and countries initiated vaccine development promptly. The Pfizer-BioNTech vaccine was the first vaccine approved for use in Canada, and distribution began in late December of 2020 in that country [3]. Vaccine production is complex, worldwide demand is extreme, and distribution logistics and vaccine storage is non-trivial. As a result, it will take many months before everyone receives a vaccine. In Canada, the priority recipients of the vaccine are generally the elderly and healthcare workers, and individual provinces determine the exact distribution order [4]. In the province of Ontario for example, health care workers in hospitals, long-term care homes, and retirement homes were amongst the first to receive vaccination [5].

Business Problem

Data scientists can play a valuable role in understanding the geographic and demographic distributions of COVID-19 infections at a detailed level. Data can be used to show whether COVID-19 infections are more prevalent in certain areas or populations rather than others.

I examine COVID-19 infections in the city of Toronto, the capital of the Canadian province of Ontario. I am interested in the rates of COVID-19 infection and I will sort these data by geographic area within the city to see if certain neighbourhoods have more infections than others, or if the distribution is random. I will generate choropleth maps to categorize the distribution of COVID-19 cases. To enhance this analysis, I will query Foursquare for neighbourhoods of interest to see the types and abundances of venues in these neighbourhoods, and determine if there are any correlations between the types of venues and the COVID-19 infection rates.

Secondly, I will attempt to understand the different regions of Toronto in terms of demographic characteristics. I will compare COVID-19 infection rates to a variety of census data, including population statistics, household income, and education level. If there are regions of the city that have particularly high or low COVID-19 cases, I will attempt to identify any demographic factors of significance in these areas. I will use an assortment of statistical and mapping analysis techniques to determine relevance of demographic factors.

The target audience for the study is government officials responsible for public health measures designed to limit the spread of COVID-19 and to coordinate vaccination distribution. My hope is that trends will become apparent, such as particular geographic

regions, common venues, or demographic factors correlate with high COVID-19 case rates. If I am successful in identifying factors that correlate with high COVID-19 infection rates, then these correlations could be used to help prioritize neighbourhoods for health measures and vaccine distributions.

Data

The city of Toronto maintains a strong online database of city-specific information. I will reference three key databases from the city of Toronto and augment with data obtained from Foursquare. The main data sources are listed in Table 1.

| Item | Description | Location |
|-------------------------------------------|-------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| COVID-19 data for city of Toronto | Daily update of COVID-19 cases in Toronto broken down by neighbourhood in Excel and PDF formats | https://www.toronto.ca/home/covid-19/covid-19-latest-city-of-toronto-news/covid-19-status-of-cases-in-toronto/ |
| City of Toronto neighbourhood profiles | Toronto census data in CSV format | https://open.toronto.ca/dataset/neighbourhood-profiles/ |
| Toronto neighbourhood spatial information | Geojson file of 140 Toronto neighbourhoods | https://open.toronto.ca/dataset/neighbourhoods/ |
| Foursquare queries | Query results of venue information for specific neighbourhoods | Results returned from Foursquare, converted to dataframe in Python |

Table 1. Key data sources utilized in project.

COVID-19 infection data for Toronto is updated daily. The COVID-19 data I employ was retrieved on January 15, 2021. The dataset consists of neighbourhood name and ID number, cumulative case count per neighbourhood, and cumulative infection rate per 100,000 people per neighbourhood. An example is shown in Figure 1.

| Neighbourhood ID | Neighbourhood Name | Rate per 100,000 people | Case Count |
|------------------|----------------------------------|-------------------------|------------|
| 138 | Eglinton East | 3464.172813 | 789 |
| 47 | Don Valley Village | 1829.876899 | 495 |
| 38 | Lansing-Westgate | 1336.302895 | 216 |
| 9 | Edenbridge-Humber Valley | 2471.837786 | 384 |
| 44 | Flemington Park | 3720.421283 | 816 |
| 59 | Danforth East York | 1222.351572 | 210 |
| 129 | Agincourt North | 2109.02346 | 614 |
| 99 | Mount Pleasant East | 846.4977645 | 142 |
| 137 | Woburn | 3791.717304 | 2028 |
| 102 | Forest Hill North | 2045.915977 | 262 |
| 111 | Rockcliffe-Smythe | 3205.070574 | 713 |
| 130 | Milliken | 2009.634201 | 534 |
| 55 | Thorncliffe Park | 4368.012128 | 922 |
| 7 | Willowridge-Martingrove-Richview | 2735.150749 | 606 |

Figure 1. Sample of COVID-19 data for neighbourhoods in the city of Toronto.

The Toronto neighbourhood profile dataset is an extensive collection of demographic data collected as part of a 2016 census. This is the most recent census data available for the city. One assumption of the project is that demographic trends have not changed in a significant manner for Toronto neighbourhoods since 2016. An example of the dataset is shown in Figure 2. The dataset may be indexed by either Neighbourhood ID or name, which makes it possible to compare with the COVID-19 data. The data contains over 2000 demographic classifications, and not all will be considered in this study. Broad categories of classification include education level, household income, age and gender, housing density, mother tongue, method of transport to employment, and citizenship.

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T | U | V | W | X | Y | Z | | | | |
|----|---------------|---------------------|------------------------------|-----------------|-----------|---------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------|-------------|-----------|-----------|------|------|------|
| 1 | Category | Topic | Data Source: Characteristic | City of Toronto | Age/Count | No | Age/Count | No | Alderwood | Annex | Bankury | Barrett | Bathurst | Bay | Street | C | Bayview | Villa | Bedford | Park | Beechwood | Bendale | Birchfield | Black Creek | Blair | | | | |
| 2 | Neighbourhood | Information | City of Toronto | Neighbourhood | Number | 329 | 328 | 20 | 95 | 42 | 34 | 76 | 52 | 49 | 39 | 132 | 127 | 122 | 121 | 120 | 119 | 118 | 117 | 116 | 115 | 114 | | | |
| 3 | Population | and dwellings | Census Profil Population | City of Toronto | 2,615,000 | 301,797 | 13,085 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | 2,615,000 | | | |
| 4 | Population | and dwellings | Census Profil Population | City of Toronto | 8,000 | 1,300 | 4,400 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | 8,000 | | | |
| 5 | Population | and dwellings | Census Profil Population | City of Toronto | 3,506 | 3,806 | 3,806 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | 3,506 | | | |
| 6 | Population | and dwellings | Census Profil Total private | 1,179,057 | 9,371 | 8,535 | 4,732 | 18,109 | 12,473 | 6,418 | 18,436 | 10,111 | 9,532 | 4,698 | 8,607 | 2,650 | 10,766 | 9,198 | 7,671 | 9,198 | 7,324 | 7,324 | 7,324 | 7,324 | 7,324 | 7,324 | | | |
| 7 | Population | and dwellings | Census Profil Private dwl | 1,112,929 | 9,120 | 8,136 | 4,616 | 15,934 | 12,144 | 6,089 | 15,074 | 9,532 | 4,698 | 8,607 | 2,650 | 10,766 | 9,198 | 7,671 | 9,198 | 7,324 | 7,324 | 7,324 | 7,324 | 7,324 | 7,324 | | | | |
| 8 | Population | and dwellings | Census Profil Population d | 4,334 | 3,929 | 3,034 | 2,435 | 10,863 | 2,775 | 3,377 | 14,937 | 4,195 | 3,240 | 3,614 | 4,031 | 3,765 | 3,765 | 3,765 | 3,765 | 3,765 | 3,765 | 3,765 | 3,765 | 3,765 | 3,765 | | | | |
| 9 | Population | and dwellings | Census Profil Private dwl | 4,082 | 7,41 | 4,871 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | 2,481 | | | | |
| 10 | Population | Age characteristics | Census Profil Children (0-1) | 398,135 | 3,840 | 3,075 | 1,760 | 2,860 | 3,605 | 2,335 | 1,695 | 2,415 | 1,515 | 4,555 | 1,120 | 4,550 | 3,345 | 4,600 | 3,345 | 4,600 | 3,345 | 4,600 | 3,345 | 4,600 | 3,345 | 4,600 | | | |
| 11 | Population | Age characteristics | Census Profil Youth (15-24) | 269,107 | 3,705 | 3,360 | 1,238 | 3,750 | 2,730 | 1,948 | 6,860 | 2,505 | 1,635 | 3,210 | 855 | 4,605 | 2,440 | 3,290 | 3,290 | 3,290 | 3,290 | 3,290 | 3,290 | 3,290 | 3,290 | 3,290 | | | |
| 12 | Population | Age characteristics | Census Profil Working Age | 1,229,555 | 11,305 | 9,965 | 5,220 | 15,040 | 10,810 | 6,655 | 13,065 | 10,310 | 4,490 | 8,410 | 2,750 | 12,050 | 9,075 | 8,525 | 9,075 | 8,525 | 9,075 | 8,525 | 9,075 | 8,525 | 9,075 | 8,525 | | | |
| 13 | Population | Age characteristics | Census Profil Pre-retirement | 336,670 | 4,230 | 3,265 | 1,822 | 3,480 | 3,555 | 2,030 | 1,760 | 2,540 | 1,825 | 3,075 | 885 | 3,535 | 3,520 | 2,425 | 3,520 | 2,425 | 3,520 | 2,425 | 3,520 | 2,425 | 3,520 | 2,425 | | | |
| 14 | Population | Age characteristics | Census Profil Senior | 426,651 | 6,045 | 4,105 | 2,015 | 5,510 | 5,770 | 2,420 | 2,420 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | 3,215 | | | | |
| 15 | Population | Age characteristics | Census Profil Older Senior | 66,000 | 925 | 555 | 320 | 1,040 | 1,640 | 710 | 330 | 610 | 740 | 660 | 145 | 900 | 630 | 370 | 630 | 370 | 630 | 370 | 630 | 370 | 630 | 370 | | | |
| 16 | Population | Age characteristics | Census Profil 0 to 04 | 69,895 | 665 | 575 | 360 | 570 | 435 | 470 | 455 | 675 | 180 | 570 | 570 | 570 | 570 | 570 | 570 | 570 | 570 | 570 | 570 | 570 | 570 | 570 | | | |
| 17 | Population | Age characteristics | Census Profil Male: 05 to C | 69,350 | 695 | 540 | 270 | 365 | 660 | 355 | 230 | 395 | 260 | 795 | 195 | 790 | 580 | 800 | 800 | 800 | 800 | 800 | 800 | 800 | 800 | 800 | | | |
| 18 | Population | Age characteristics | Census Profil Male: 15 to 19 | 64,945 | 660 | 460 | 225 | 570 | 650 | 415 | 415 | 415 | 415 | 415 | 260 | 880 | 120 | 775 | 585 | 775 | 585 | 775 | 585 | 775 | 585 | 775 | 585 | | |
| 19 | Population | Age characteristics | Census Profil Male: 20 to 24 | 74,740 | 840 | 780 | 285 | 465 | 715 | 490 | 585 | 320 | 385 | 880 | 120 | 655 | 585 | 655 | 655 | 655 | 655 | 655 | 655 | 655 | 655 | 655 | | | |
| 20 | Population | Age characteristics | Census Profil Male: 25 to 29 | 97,415 | 1015 | 1000 | 305 | 125 | 125 | 700 | 530 | 2485 | 735 | 445 | 765 | 215 | 1370 | 650 | 840 | 840 | 840 | 840 | 840 | 840 | 840 | 840 | | | |
| 21 | Population | Age characteristics | Census Profil Male: 30 to 34 | 113,905 | 1015 | 1045 | 355 | 280 | 645 | 465 | 215 | 1075 | 405 | 405 | 405 | 225 | 1135 | 530 | 725 | 725 | 725 | 725 | 725 | 725 | 725 | 725 | | | |
| 22 | Population | Age characteristics | Census Profil Male: 35 to 39 | 108,951 | 835 | 820 | 410 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | 1610 | | | |
| 23 | Population | Age characteristics | Census Profil Male: 40 to 44 | 94,625 | 605 | 525 | 245 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | 1005 | | | |
| 24 | Population | Age characteristics | Census Profil Male: 45 to 49 | 86,535 | 760 | 610 | 210 | 420 | 835 | 815 | 435 | 560 | 685 | 310 | 665 | 170 | 840 | 740 | 585 | 585 | 585 | 585 | 585 | 585 | 585 | 585 | | | |
| 25 | Population | Age characteristics | Census Profil Male: 50 to 54 | 90,860 | 895 | 760 | 440 | 570 | 500 | 1010 | 535 | 500 | 605 | 390 | 780 | 225 | 950 | 835 | 660 | 660 | 660 | 660 | 660 | 660 | 660 | 660 | | | |
| 26 | Population | Age characteristics | Census Profil Male: 55 to 59 | 98,735 | 1160 | 960 | 550 | 515 | 920 | 1110 | 605 | 540 | 620 | 430 | 845 | 220 | 1010 | 1020 | 680 | 680 | 680 | 680 | 680 | 680 | 680 | 680 | | | |
| 27 | Population | Age characteristics | Census Profil Male: 60 to 64 | 88,145 | 1060 | 850 | 460 | 580 | 580 | 895 | 930 | 565 | 475 | 560 | 385 | 750 | 225 | 750 | 945 | 945 | 945 | 945 | 945 | 945 | 945 | 945 | | | |
| 28 | Population | Age characteristics | Census Profil Male: 65 to 69 | 72,740 | 925 | 710 | 390 | 755 | 730 | 385 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | 730 | | | |
| 29 | Population | Age characteristics | Census Profil Male: 70 to 74 | 60,360 | 925 | 630 | 300 | 780 | 715 | 380 | 295 | 475 | 405 | 560 | 225 | 345 | 315 | 400 | 125 | 495 | 420 | 285 | 285 | 285 | 285 | 285 | 285 | | |
| 30 | Population | Age characteristics | Census Profil Male: 75 to 79 | 42,320 | 590 | 425 | 205 | 640 | 570 | 405 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | 225 | | |
| 31 | Population | Age characteristics | Census Profil Male: 80 to 84 | 32,730 | 490 | 375 | 155 | 485 | 485 | 505 | 205 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | 200 | | |
| 32 | Population | Age characteristics | Census Profil Male: 85 to 89 | 29,475 | 395 | 345 | 125 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | 365 | |
| 33 | Population | Age characteristics | Census Profil Male: 90 to E | 25,670 | 380 | 240 | 105 | 335 | 435 | 165 | 150 | 230 | 255 | 250 | 50 | 35 | 75 | 200 | 355 | 200 | 355 | 200 | 355 | 200 | 355 | 200 | 355 | 200 | |
| 34 | Population | Age characteristics | Census Profil Male: 85 to 1 | 15,665 | 210 | 155 | 65 | 230 | 340 | 110 | 80 | 150 | 165 | 145 | 35 | 75 | 10 | 80 | 55 | 40 | 55 | 40 | 55 | 40 | 55 | 40 | 55 | 40 | |
| 35 | Population | Age characteristics | Census Profil Male: 90 to F | 6,185 | 100 | 45 | 35 | 80 | 175 | 105 | 35 | 75 | 75 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 36 | Population | Age characteristics | Census Profil Male: 95 to 1 | 1,280 | 25 | 20 | 10 | 20 | 20 | 20 | 20 | 20 | 20 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 37 | Population | Age characteristics | Census Profil Female: 1 | 112,930 | 925 | 8135 | 4620 | 15935 | 6086 | 15075 | 9353 | 4700 | 8610 | 2655 | 10765 | 9205 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 | 7320 |
| 38 | Population | Age characteristics | Census Profil Female: 2 | 105,670 | 1290 | 975 | 515 | 935 | 1175 | 630 | 665 | 750 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | 520 | | |
| 39 | Population | Age characteristics | Census Profil Female: 3 | 94,660 | 1160 | 915 | 485 | 485 | 1005 | 1005 | 600 | 600 | 365 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | 235 | | |
| 40 | Population | Age characteristics | Census Profil Female: 4 | 81,600 | 1070 | 795 | 400 | 940 | 940 | 475 | 470 | 425 | 660 | 660 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | 455 | |
| 41 | Population | Age characteristics | Census Profil Female: 5 | 70,970 | 985 | 689 | 310 | 650 | 650 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | 305 | | |
| 42 | Population | Age characteristics | Census Profil Female: 6 | 51,288 | 690 | 480 | 210 | 700 | 700 | 275 | 320 | 320 | 400 | 400 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | 505 | |
| 43 | Population | Age characteristics | Census Profil Female: 7 | 43,430 | 575 | 405 | 180 | 565 | 730 | 280 | 250 | 355 | 400 | 400 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | 90 | |
| 44 | Population | Age characteristics | Census Profil Female: 8 | 34,965 | 485 | 350 | 210 | 425 | 650 | 285 | 170 | 315 | 330 | 305 | 60 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | 50 | |
| 45 | Population | Age characteristics | Census Profil Female: 9 | 35,135 | 350 | 205 | 130 | 205 | 305 | 165 | | | | | | | | | | | | | | | | | | | |

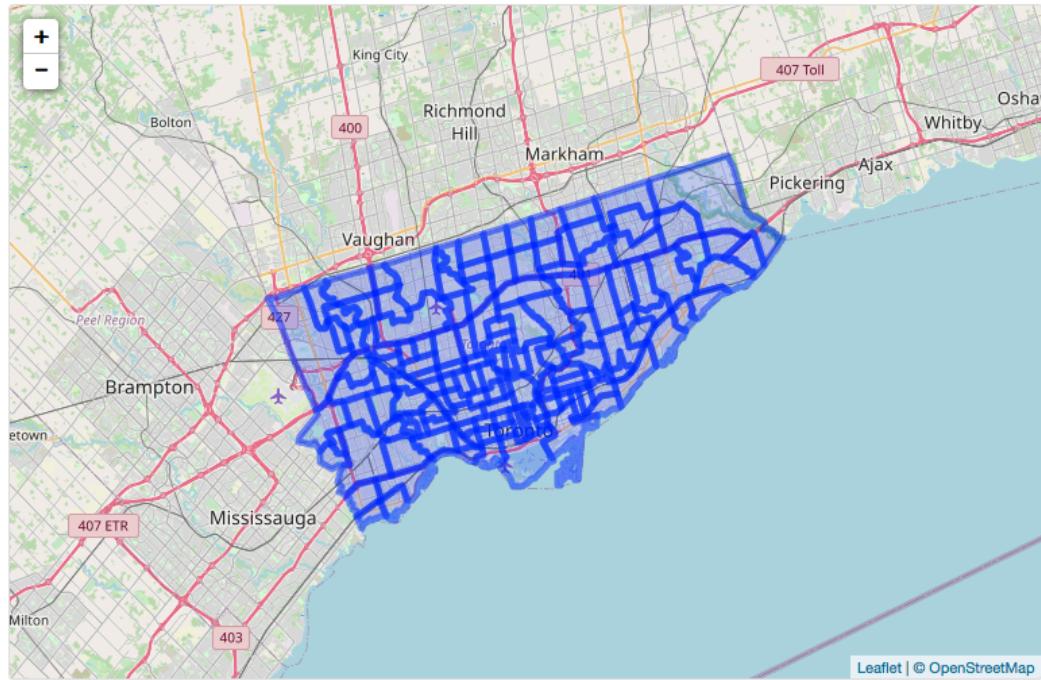


Figure 3. Map view of Geojson file used in this study.

Once I identify neighbourhoods of high or low rates of COVID-19 infection, I will run queries in Foursquare to understand the types and abundances of venues in each neighbourhood. Figure 4 is a sample of a Foursquare result converted to a data frame for analysis.

| [62] : | Cluster Labels | Neighborhood | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue | 5th Most Common Venue | 6th Most Common Venue | 7th Most Common Venue |
|----------|----------------|-----------------------------------|-----------------------|------------------------|----------------------------|------------------------|------------------------|------------------------|-----------------------|
| 0 | 1 | Black Creek | Home Service | Hotel | Fast Food Restaurant | Intersection | Electronics Store | Coffee Shop | Italian Restaurant |
| 1 | 1 | Forest Hill South | Pizza Place | Trail | Gift Shop | Intersection | Israeli Restaurant | Coffee Shop | Korean Restaurant |
| 2 | 4 | Humbermede | Golf Course | Furniture / Home Store | Coffee Shop | Department Store | Dance Studio | Gastropub | Food & Drink Shop |
| 3 | 3 | Maple Leaf | Park | Trail | Convenience Store | Furniture / Home Store | Food & Drink Shop | Fish & Chips Shop | Fast Food Restaurant |
| 4 | 0 | Mount Olive-Silverstone-Jamestown | Park | Baseball Field | Caribbean Restaurant | Trail | Dance Studio | Furniture / Home Store | Food & Drink Shop |
| 5 | 1 | Rosedale-Moore Park | Coffee Shop | Fish & Chips Shop | Breakfast Spot | Pizza Place | Fast Food Restaurant | Intersection | Trail |
| 6 | 1 | Runnymede-Bloor West Village | Asian Restaurant | Café | Food & Drink Shop | Diner | Market | Coffee Shop | Sushi Restaurant |
| 7 | 1 | The Beaches | Fast Food Restaurant | Italian Restaurant | Ice Cream Shop | Burrito Place | Liquor Store | Food & Drink Shop | Fish & Chips Shop |
| 8 | 2 | Thistletown-Beaumont Heights | Convenience Store | Home Service | Construction & Landscaping | Dance Studio | Furniture / Home Store | Food & Drink Shop | Fish & Chips Shop |
| 9 | 1 | Woodbine Corridor | Fast Food Restaurant | Italian Restaurant | Ice Cream Shop | Burrito Place | Liquor Store | Food & Drink Shop | Fish & Chips Shop |

Figure 4. Example of data frame created from Foursquare data query.

Methodology

Three datasets – COVID-19 data, demographic data, and a Geojson file - were downloaded from the City of Toronto website. The COVID-19 and demographic data each contained columns of neighbourhood name and neighbourhood number (1-140), so merging the datasets was possible. The Geojson file referenced the same neighbourhood number, allowing the statistical data to be referenced to the spatial data. Data cleaning and manipulation were completed in both Excel and Python. For many of the figures included in this report, abbreviations for data categories are employed and will be explained as required. The most important is ‘Covid19Rate,’ representing the number of confirmed COVID-19 infections per 100,000 residents of each neighbourhood. Covid19CaseCount is the total cumulative cases of COVID-19 for the given neighbourhood.

The most fundamental question to answer is if certain neighbourhoods have higher rates of COVID-19 infections than others. Figures 5 and 6 show the neighbourhoods with the lowest and highest number of COVID-19 infections per 100,000 residents in table form. There are significant differences, as the neighbourhood with the highest infection rate (Thistletown-Beaumont Heights) has a rate of more than ten times the neighbourhood with the lowest infection rate (The Beaches).

| Covid19Rate | |
|-----------------------------------|------------|
| Neighbourhood | |
| Bridle Path-Sunnybrook-York Mills | 938.916469 |
| Newtonbrook East | 938.062993 |
| Lawrence Park South | 935.502998 |
| Willowdale East | 854.582226 |
| Mount Pleasant East | 846.497764 |
| Rosedale-Moore Park | 831.620704 |
| Woodbine Corridor | 813.332270 |
| Forest Hill South | 810.659709 |
| Runnymede-Bloor West Village | 645.481629 |
| The Beaches | 625.956322 |

Figure 5. Neighbourhoods with lowest rates of COVID-19 infection.

| Covid19Rate | |
|-----------------------------------|-------------|
| Neighbourhood | |
| Thistletown-Beaumont Heights | 6650.579151 |
| Mount Olive-Silverstone-Jamestown | 6463.555259 |
| Maple Leaf | 6201.167046 |
| Black Creek | 6095.597369 |
| Humbermede | 5969.765198 |
| West Humber-Clairville | 5586.575408 |
| Humber Summit | 5565.399485 |
| Glenfield-Jane Heights | 5559.017415 |
| Weston | 5491.329480 |
| Downsview-Roding-CFB | 5232.226406 |

Figure 6. Neighbourhoods with highest rates of COVID-19 infection.

Next Foursquare was used to understand the number and variety of venues present within each neighbourhood of interest. Foursquare queries were analyzed for the five neighbourhoods with the highest infection rates and the five neighbourhoods with the lowest infection rates. The number of venues per neighbourhood, and the types and abundance were considered. It was difficult to identify clear trends from the data. This will be expanded upon in the Results section.

The demographic data were examined for factors correlating with COVID-19 infection rates. A subset of the full demographic suite was chosen for investigation, specifically data related to visible minority status, age, household income, employment, education, and commute to work.

For each demographic category, a similar workflow was employed. A correlation matrix was generated and plotted in tables and as heatmaps. The critical information is the correlation between Covid19Rate and the demographic feature of interest. Regression plots were also examined to visualize correlations that warrant further attention. For brevity, not all figures are captured here.

The correlation matrices for the visible minority data were split for ease of display and shown in figures 7 and 8. A correlation of 1.0 represents a perfect positive correlation, where the y-axis variable increases as the x-axis variable increases. Conversely, a value of -1.0 means the y-axis variable decreases at a proportional rate as the x-axis increases. A value of zero represents no correlation. Of primary interest is the row Covid19Rate, showing the correlation between the infection rate and the visible minority populations.

| | NeighbourhoodNumber | Covid19Rate | Covid19CaseCount | VMAboriginal | VMLatinAmerica | VMBLack | VMArab | VMNotSpecified | VMMultiple | VMNot |
|---------------------|---------------------|-------------|------------------|--------------|----------------|-----------|-----------|----------------|------------|-----------|
| NeighbourhoodNumber | 1.000000 | -0.160640 | -0.027007 | 0.260068 | -0.221717 | 0.022347 | -0.070637 | 0.083444 | 0.183580 | -0.073978 |
| Covid19Rate | -0.160640 | 1.000000 | 0.730923 | -0.096396 | 0.533566 | 0.654493 | 0.166010 | 0.493298 | 0.311711 | -0.334786 |
| Covid19CaseCount | -0.027007 | 0.730923 | 1.000000 | 0.223191 | 0.602199 | 0.851278 | 0.462577 | 0.778155 | 0.776954 | 0.040880 |
| VMAboriginal | 0.260068 | -0.096396 | 0.223191 | 1.000000 | 0.173971 | 0.322458 | 0.206579 | 0.381450 | 0.435464 | 0.530763 |
| VMLatinAmerica | -0.221717 | 0.533566 | 0.602199 | 0.173971 | 1.000000 | 0.649863 | 0.198034 | 0.431688 | 0.488615 | 0.192868 |
| VMBLack | 0.022347 | 0.654493 | 0.851278 | 0.322458 | 0.649863 | 1.000000 | 0.395387 | 0.806073 | 0.736919 | -0.062478 |
| VMArab | -0.070637 | 0.166010 | 0.462577 | 0.206579 | 0.198034 | 0.395387 | 1.000000 | 0.364912 | 0.524611 | 0.165470 |
| VMNotSpecified | 0.083444 | 0.493298 | 0.778155 | 0.381450 | 0.431688 | 0.806073 | 0.364912 | 1.000000 | 0.725973 | -0.009871 |
| VMMultiple | 0.183580 | 0.311711 | 0.776954 | 0.435464 | 0.488615 | 0.736919 | 0.524611 | 0.725973 | 1.000000 | 0.229005 |
| VMNot | -0.073978 | -0.334786 | 0.040880 | 0.530763 | 0.192868 | -0.062478 | 0.165470 | -0.009871 | 0.229005 | 1.000000 |

Figure 7. Correlation matrix for a portion of visible minority demographic data. The categories shown are visible minority Aboriginal, Latin American, Black, Arab, unspecified, multiple visible minorities, or not a visible minority.

| | NeighbourhoodNumber | Covid19Rate | Covid19CaseCount | VMSouthAsian | VMChinese | VMFilipino | VMSoutheastAsian | VMWestAsian | VMKorean | VMJapanese |
|---------------------|---------------------|-------------|------------------|--------------|-----------|------------|------------------|-------------|-----------|------------|
| NeighbourhoodNumber | 1.000000 | -0.160640 | -0.027007 | 0.238937 | 0.218038 | 0.172835 | -0.171758 | -0.120980 | -0.179417 | 0.087104 |
| Covid19Rate | -0.160640 | 1.000000 | 0.730923 | 0.377711 | -0.194829 | 0.410632 | 0.458677 | -0.013907 | -0.216567 | -0.423823 |
| Covid19CaseCount | -0.027007 | 0.730923 | 1.000000 | 0.740276 | 0.083278 | 0.666958 | 0.519743 | 0.214637 | -0.015151 | -0.060593 |
| VMSouthAsian | 0.238937 | 0.377711 | 0.740276 | 1.000000 | 0.182574 | 0.537870 | 0.059744 | 0.278101 | -0.021373 | 0.025877 |
| VMChinese | 0.218038 | -0.194829 | 0.083278 | 0.182574 | 1.000000 | 0.083815 | -0.014296 | 0.413517 | 0.408836 | 0.410311 |
| VMFilipino | 0.172835 | 0.410632 | 0.666958 | 0.537870 | 0.083815 | 1.000000 | 0.159245 | 0.234661 | 0.070573 | -0.028514 |
| VMSoutheastAsian | -0.171758 | 0.458677 | 0.519743 | 0.059744 | -0.014296 | 0.159245 | 1.000000 | 0.024845 | -0.035139 | -0.044588 |
| VMWestAsian | -0.120980 | -0.013907 | 0.214637 | 0.278101 | 0.413517 | 0.234661 | 0.024845 | 1.000000 | 0.834930 | 0.494523 |
| VMKorean | -0.179417 | -0.216567 | -0.015151 | -0.021373 | 0.408836 | 0.070573 | -0.035139 | 0.834930 | 1.000000 | 0.651388 |
| VMJapanese | 0.087104 | -0.423823 | -0.060593 | 0.025877 | 0.410311 | -0.028514 | -0.044588 | 0.494523 | 0.651388 | 1.000000 |

Figure 8. Correlation matrix for a second portion of visible minority demographic data. The categories shown are visible minority South Asian, Chinese, Filipino, Southeast Asian, West Asian, Korean, and Japanese.

The visible minorities with the highest correlation to COVID-19 infections per 100,000 are Black (0.65), and Latin American (0.53). Japanese has the highest negative correlation at -0.45. Regression plots for visible minority Black and visible minority Japanese are shown in figures 9 and 10.

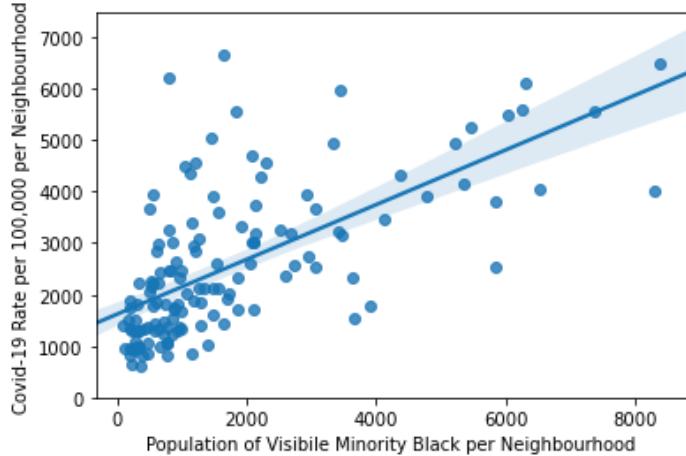


Figure 9. Regression plot of COVID-19 cases per 100,000 people for visible minority Black populations showing a strong positive correlation.

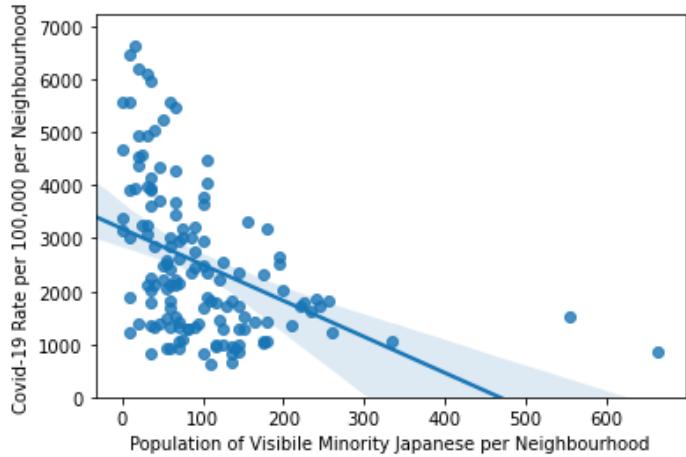


Figure 10. Regression plot of COVID-19 cases per 100,000 people for visible minority Japanese populations showing a negative correlation.

Total household income was separated into bins. The heat map of the correlation matrix for total household income and COVID-19 infections is shown in Figure 11. The bins chosen for household income are below \$25,000, \$25,000-\$40,000, \$40,000-60,000, \$60,000-\$80,000, and greater than \$80,000. The most important row to consider is Covid19Rate. Neither of the income bins has a very high correlation with Covid19Rate, with the max value 0.20 for \$25,000-\$40,000. But the highest income bin (\$80,000 plus) has the largest negative correlation at -0.26.

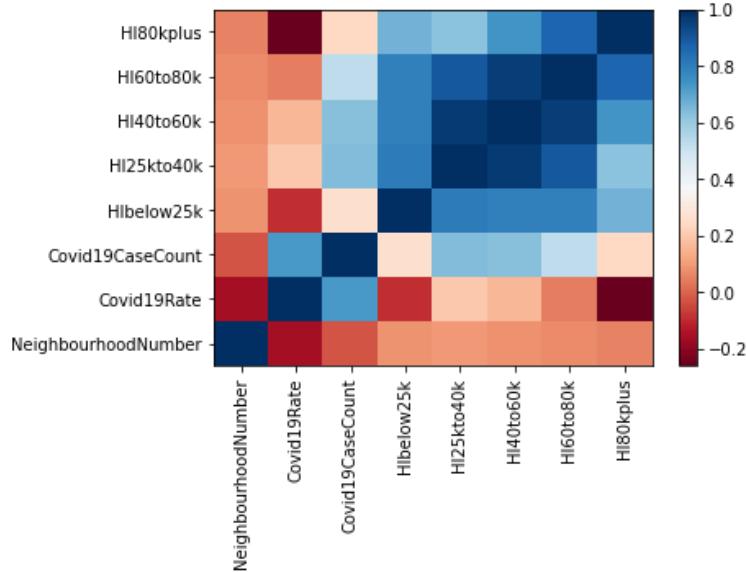


Figure 11. Heatmap of correlation matrix for total household income data.

Age of residents was grouped in bins of 0-14 years, 15-25 years, 25-54 years, 55-64 years, 65 plus, and 85 plus, the latter two subject to some overlap due to lack of specific differentiation. The heat map of the correlation matrix for age and COVID-19 infections is shown in Figure 12. Correlations do not differ greatly by age bin, but the highest positive correlations are for the two youngest bins, representing the range of infants, to school age, to young professionals.

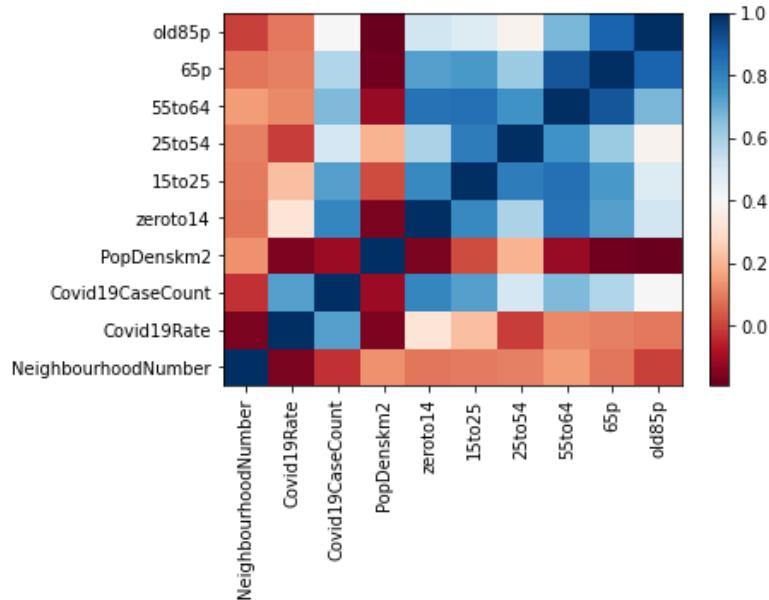


Figure 12. Heatmap of correlation matrix for age data.

An assortment of data describing commuting to work was analyzed, including duration of commute and method of transportation. There are high positive correlations between COVID-19 infection rates and both long commutes to work and commuting as a passenger

in a private vehicle, the second of which is shown in Figure 13. The strongest negative correlation is with those who cycle to work. There is no significant correlation between taking public transit and COVID-19 infections, as shown in Figure 14.

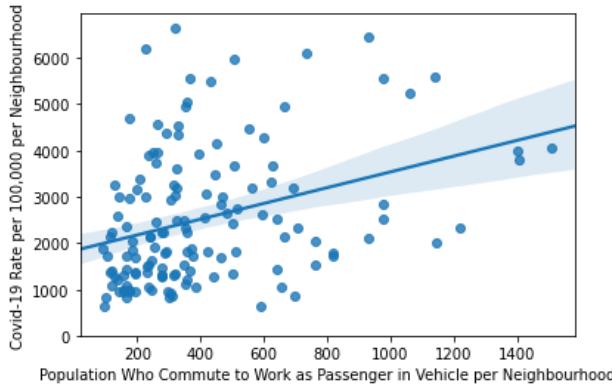


Figure 13. Regression plot showing a positive correlation between commuting to work as a passenger in a private vehicle and COVID-19 infections.

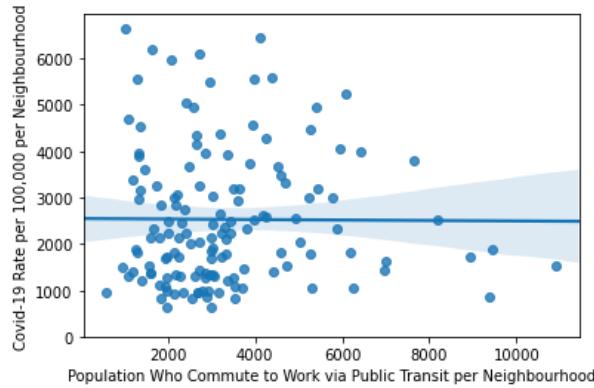


Figure 14. Regression plot showing no correlation between commuting to work via public transit and COVID-19 infections.

Employment data were also examined, and the correlation matrix is shown in table form in Figure 15. COVID-19 infection rate has a strong positive correlation with unemployment, and a strong negative correlation with employment. The trend for unemployment is shown in Figure 16.

| | NeighbourhoodNumber | Covid19Rate | Covid19CaseCount | EmploymentRate | UnemploymentRate |
|---------------------|---------------------|-------------|------------------|----------------|------------------|
| NeighbourhoodNumber | 1.000000 | -0.160640 | -0.027007 | 0.044850 | 0.091673 |
| Covid19Rate | -0.160640 | 1.000000 | 0.730923 | -0.518157 | 0.605126 |
| Covid19CaseCount | -0.027007 | 0.730923 | 1.000000 | -0.376483 | 0.486920 |
| EmploymentRate | 0.044850 | -0.518157 | -0.376483 | 1.000000 | -0.765290 |
| UnemploymentRate | 0.091673 | 0.605126 | 0.486920 | -0.765290 | 1.000000 |

Figure 15. Correlation matrix for employment and unemployment rates.

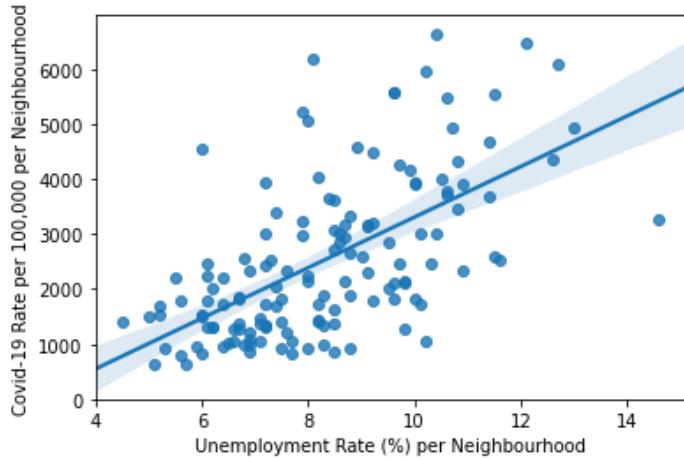


Figure 16. Regression plot of COVID-19 cases per 100,000 people vs. unemployment rate, showing a strong positive correlation.

The final demographic category considered was education level. Education was broken down into several bins, considering only adults aged 25-64. The categories and abbreviations are: no high school diploma (A25to64NoCert), high school diploma (A25to64HSdip), trade certificate or diploma (A25to64Trade), diploma from a non-university college or CEGEP (A25to64non-uniDip), university certificate or diploma below bachelor level (A25to64UniDipbelowBach), and university diploma of bachelor level or greater (A25to64UniBachorhigher). The heatmap of the correlation matrix is shown in Figure 17. Considering the row of Covid19rate, the highest positive correlations are associated with education less than high school and certified trades. University education bachelor or higher, has a negative correlation.

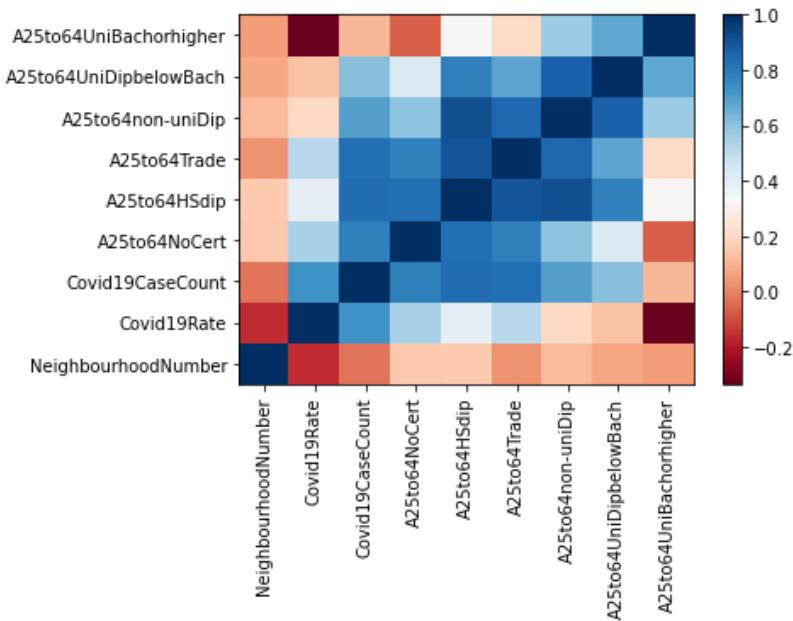


Figure 17. Heatmap of correlation matrix for education levels

Results

Neighbourhoods in Toronto with the highest and lowest number of COVID-19 infections per 100,000 people were identified in the Methodology section and shown in table form. These data are repeated in bar charts and are shown below. Figure 18 shows the neighbourhoods with the lowest number of infections per 100,000 people, while figure 19 shows the neighbourhoods with the highest number of infections. This is better visualized in the choropleth map shown in Figure 20. The highest rates are observed in neighbourhoods to the northwest, at the city limits. Case rates are also high in eastern neighbourhoods, and lowest in the centre and in the lakeside downtown core.

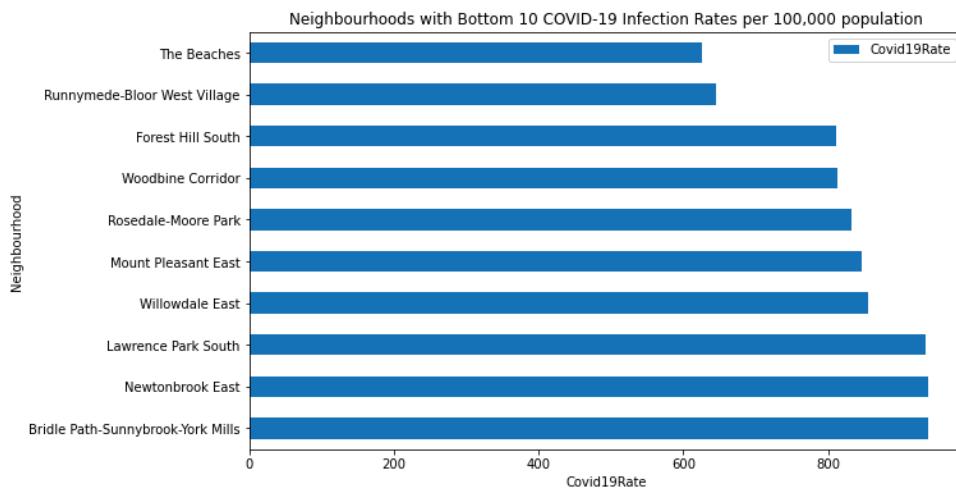


Figure 18. Neighbourhoods with lowest number of COVID-19 infections per 100,000 residents.

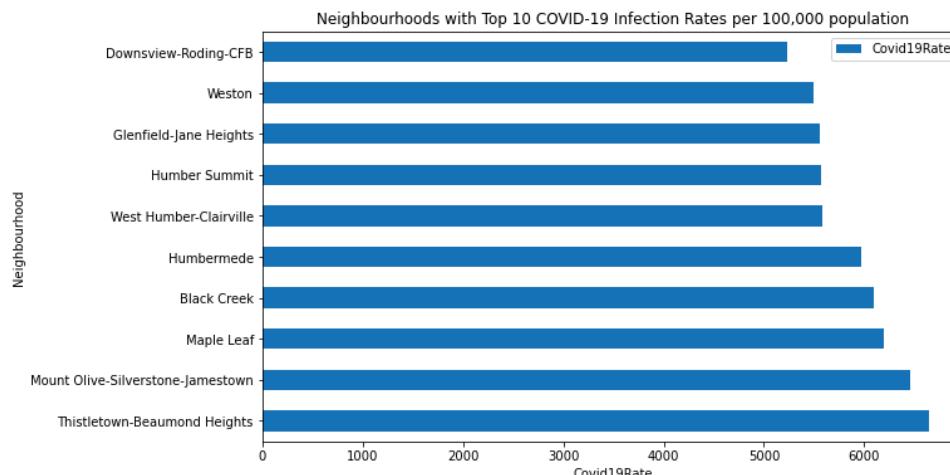


Figure 19. Neighbourhoods with highest number of COVID-19 infections per 100,000 residents.

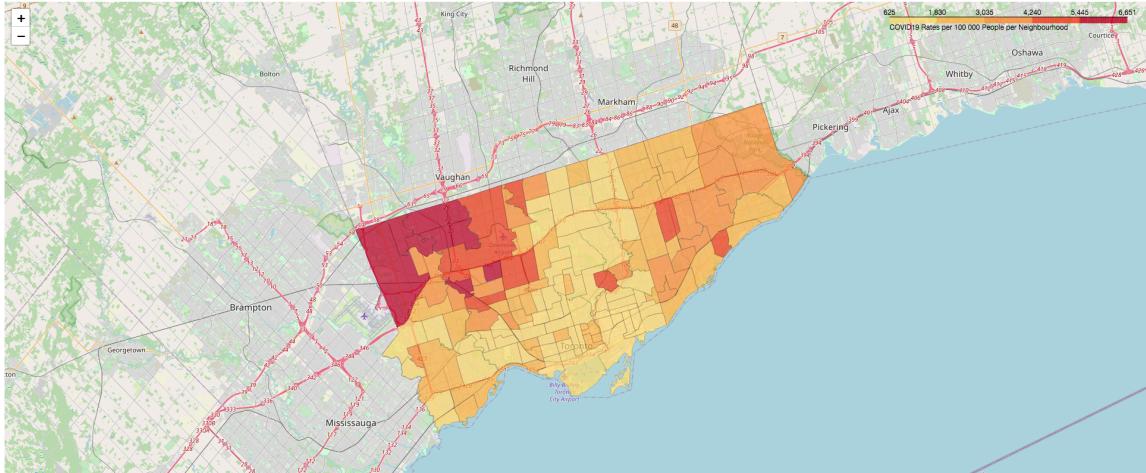


Figure 20. Choropleth map of COVID-19 infections per 100,000 residents of each Toronto neighbourhood.

The neighbourhoods with the highest and lowest rates of infection were tested with queries in Foursquare. Figure 21 shows the number of venues (of all categories) for the neighbourhoods with the top five and bottom five infections per 100,000 residents. It is apparent that the five neighbourhoods with the highest number of venues are the five neighbourhoods with the lowest rate of infection per 100,000 residents. The neighbourhoods with fewer venues have the highest infection rates.

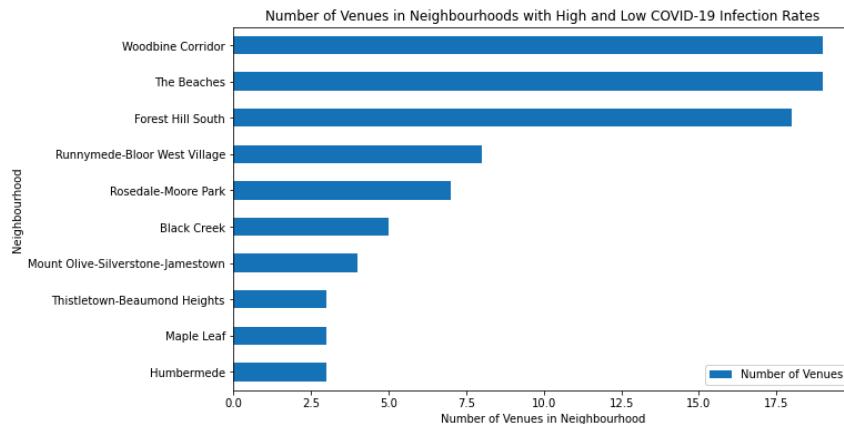


Figure 21. Number of venues (uncategorized) for the neighbourhoods with top five and bottom five COVID-19 infection rates.

Data returned from Foursquare were further examined and categorized by the top five venues per neighbourhood. This breakdown is shown in Figure 22. It is difficult to make strong classifications of the results. Four of the neighbourhoods with the lowest rates of COVID-19 infection have at least one outdoor venue in the top five venues. Runnymede-Bloor West Village is the exception. Similarly, ignoring intersection as a venue, four of the five neighbourhoods with the highest rates of COVID-19 infection have at least one outdoor venue in the top five venues. Of the indoor venues, most are an assortment of recreational facilities or food services. Since clear relationships between COVID-19 infection rates and venues returned from Foursquare could not be established, demographic data were analyzed to understand the geographic distribution of infection rates.

| Neighborhood | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue | 5th Most Common Venue |
|-----------------------------------|-----------------------|----------------------------|-----------------------|-----------------------|-----------------------|
| Black Creek | Hotel | Intersection | Fast Food Restaurant | Coffee Shop | Home Service |
| Forest Hill South | Bagel Shop | Restaurant | Pizza Place | Trail | Optical Shop |
| Humbermede | Golf Course | Coffee Shop | Department Store | Trail | Dance Studio |
| Maple Leaf | Tennis Court | Pharmacy | Park | Convenience Store | Golf Course |
| Mount Olive-Silverstone-Jamestown | Park | Baseball Field | Caribbean Restaurant | Dance Studio | Golf Course |
| Rosedale-Moore Park | Coffee Shop | Intersection | Breakfast Spot | Pizza Place | Fish & Chips Shop |
| Runnymede-Bloor West Village | Asian Restaurant | Sushi Restaurant | Food & Drink Shop | Pizza Place | Café |
| The Beaches | Fast Food Restaurant | Pizza Place | Park | Steakhouse | Sandwich Place |
| Thistletown-Beaumont Heights | Convenience Store | Construction & Landscaping | Home Service | Bank | Baseball Field |
| Woodbine Corridor | Fast Food Restaurant | Pizza Place | Park | Steakhouse | Sandwich Place |

Figure 22. Most common venue types per neighbourhood, using the subset of neighbourhoods with the top five and bottom five rates of infection.

Figures 23-27 are a series of choropleth maps constructed from the Toronto demographic data. The chosen figures help to describe the demographic characteristics associated (or not associated) with high rates of COVID-19 infections. These maps should be directly compared with Figure 20, the map of infection rates by neighbourhood. The demographic data explain COVID-19 prevalence better than the venue data obtained from Foursquare.

Figure 23 shows the population of people who identify as visible minority Black for each Toronto neighbourhood. As was suggested in the statistical analysis, neighbourhoods with high populations of people who identify as visible minority Black generally have high rates of COVID-19 infection per 100,000 residents.

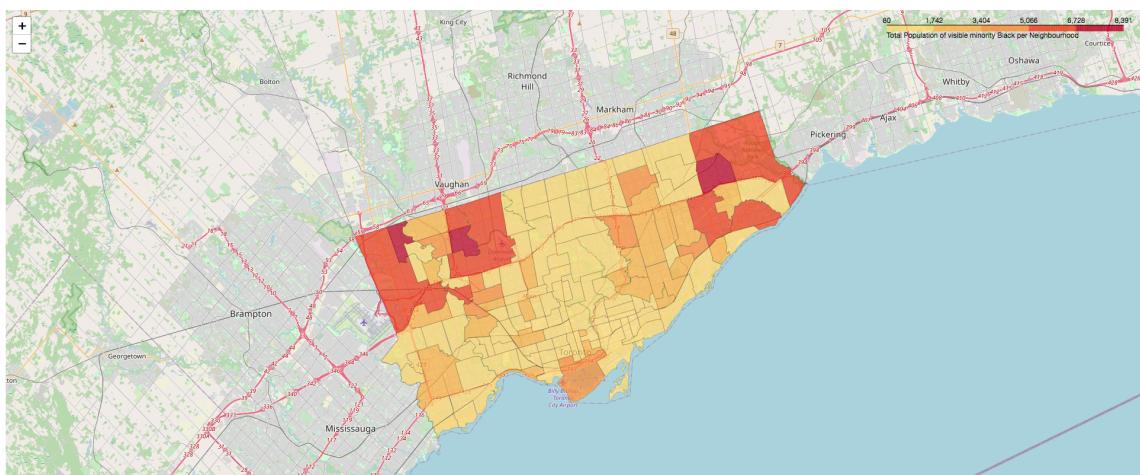


Figure 23. Choropleth map of population of people who identify as visible minority Black for each Toronto neighbourhood.

Figure 24 is a choropleth map for the population with total household income of \$80,000 or greater for each Toronto neighbourhood. Neighbourhoods with high populations of wealthy people generally do not have high rates of COVID-19 infection. The number of people per neighbourhood who commute to work as a vehicle passenger (not public transit) is presented in Figure 25. This map resembles the COVID-19 infection rate map in Figure 20. Figure 26 shows the unemployment rate by neighbourhood. Many areas with high COVID-19 infection rates (Figure 20) are also regions of high unemployment.

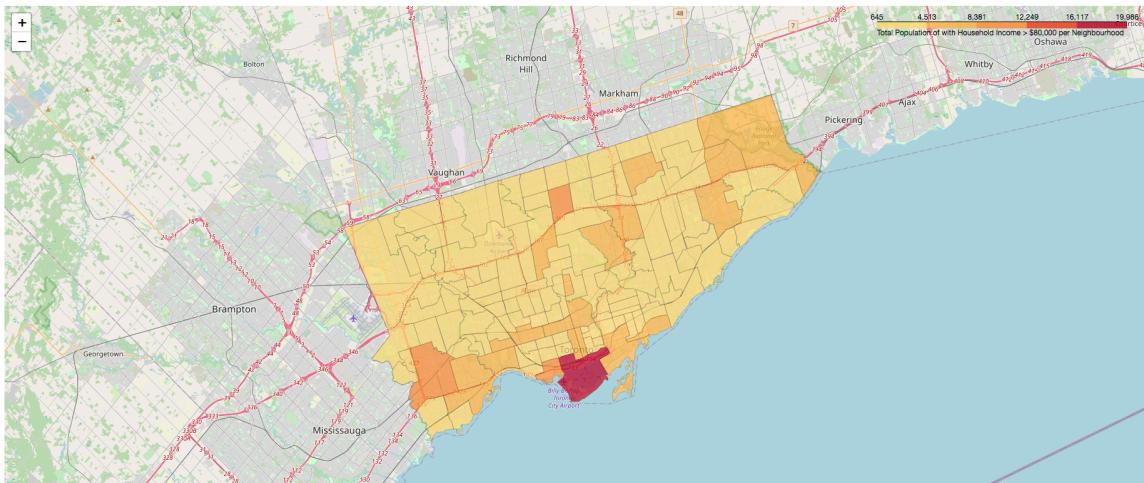


Figure 24. Choropleth map of population of people with total household income \$80,000 and greater for each Toronto neighbourhood.

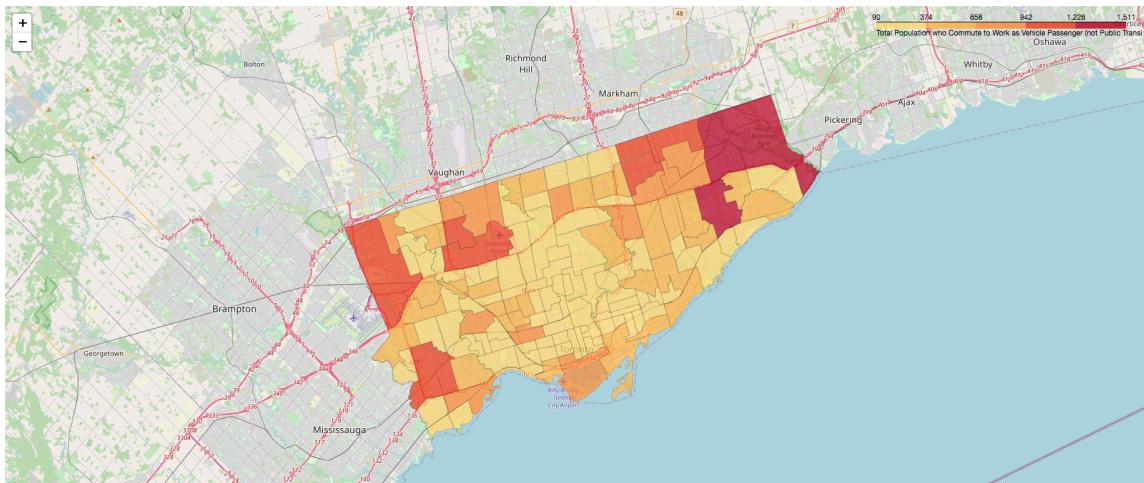


Figure 25. Choropleth map of population of people who commute to work in a private vehicle (not public transit) as a passenger for each Toronto neighbourhood.

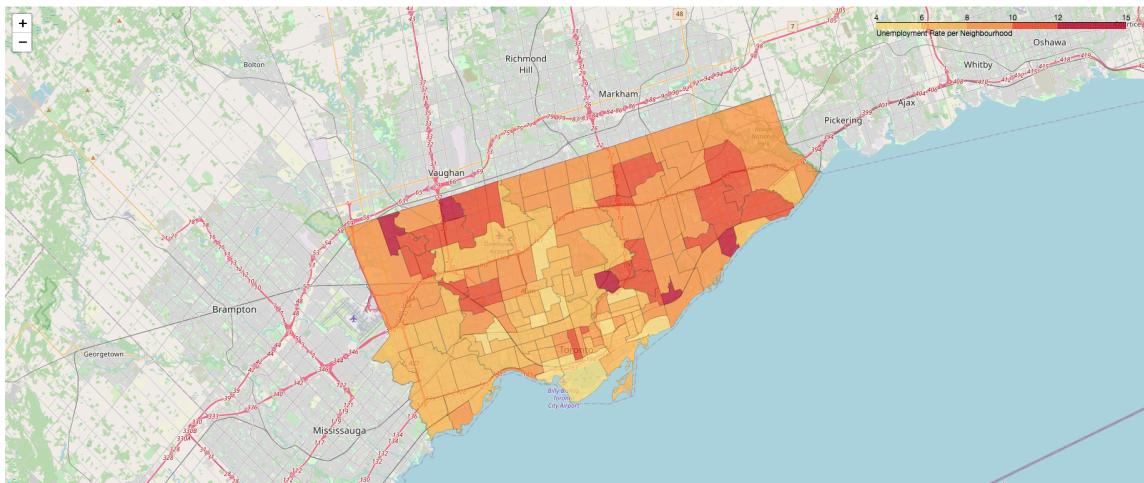


Figure 26. Choropleth map of unemployment rate for each Toronto neighbourhood.

Lastly, the number of people per neighbourhood with education level of Bachelor's Degree or higher is shown in Figure 27. The region of the most highly educated people corresponds with neighbourhoods of low rates of COVID-19 infection.

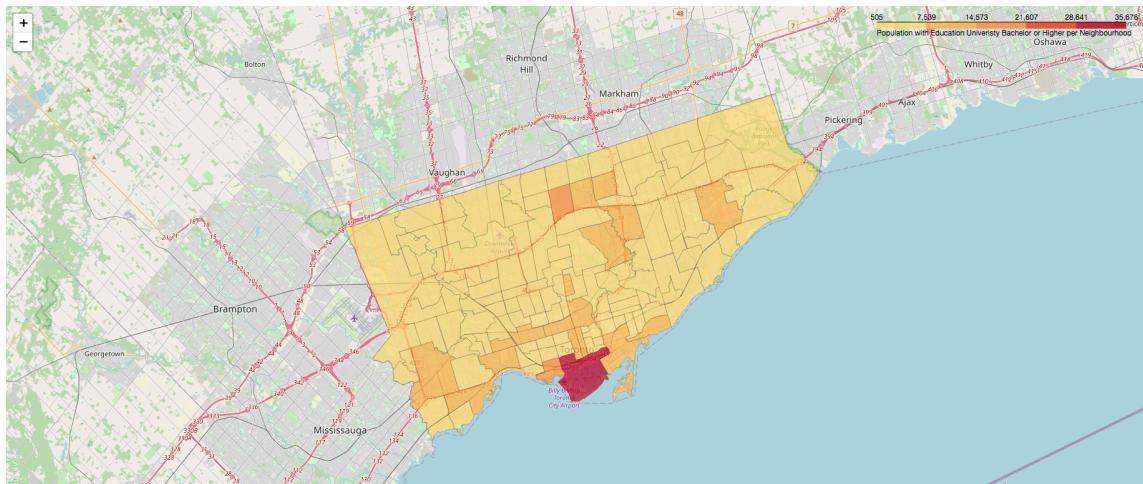


Figure 27. Choropleth map of number of people with university degrees at bachelor level or higher for each Toronto neighbourhood.

Discussion

A rich database was analyzed to understand the distribution of COVID-19 cases in Toronto. It was relatively easy to determine which neighbourhoods had recorded the highest number of COVID-19 infections. The top five neighbourhoods in terms of infection rates are Thistletown-Beaumont Heights, Mount Olive-Silverstone-Jamestown, Maple Leaf, Black Creek, and Humbermede. These have been presented in table and map form.

It was not immediately obvious why infection rates were highest in these regions. Foursquare queries of the neighbourhoods with highest and lowest rates of infections showed that neighbourhoods with the most public venues had the lowest rates of infections. It could perhaps be concluded that access to public venues is helpful in reducing COVID-19 infection rates. Perhaps access to public venues reduces gatherings in private homes for example. But the data at hand could not be used to verify this, and quarantine factors such as restricted access hours or outright closures of these venues were not clear. Further breakdowns into the nature of venues, be they indoor or outdoor in nature, did not add further clarification.

Thus demographic data were analyzed, and proved helpful in understanding the characteristics of individuals most or least prone to COVID-19 infections. Factors that strongly positively correlate with COVID-19 infections were common in neighbourhoods where infection rates are high. These factors include, but are not limited to, identifying as a visible minority black, commuting to work as a passenger in a private vehicle, and being unemployed. Conversely, factors that have strong negative correlations with infection rates were common in neighbourhoods where infection rates are low. These factors include high household income, and high education.

Conclusions

The goal of this study was to help inform policy makers responsible for public health measures. Regions with high infection rates of COVID-19 should be prioritized for measures aimed at reducing infections and for prioritized vaccinations. Based on this study, it is reasonable to conclude that neighbourhoods with the highest infection rates (Thistletown-Beaumont Heights, Mount Olive-Silverstone-Jamestown, Maple Leaf, Black Creek, and Humbermede) should be subjected to new health measures and prioritized for vaccinations.

Fortunately, the richness of the Toronto dataset allows for a more nuanced vaccination prioritization. This study suggests demographic groups that should be prioritized for vaccination and those who should wait. *Vaccination priority groups* include those who identify as visible minority Black or Latin American, occupants of households with low to moderate incomes, people aged 25 or lower, those who commute to work as a passenger or for a long duration, the unemployed, and those with low education or who work in trades. Conversely, groups that should be *low vaccination priority* include visible minority Japanese or those who are not a visible minority, occupants of households with high income, and individuals who are highly educated.

References

- [1] Wikipedia, retrieved January 21, 2021. *COVID-19 Pandemic*.
https://en.wikipedia.org/wiki/COVID-19_pandemic
- [2] World Health Organization, retrieved January 21, 2021. *WHO Coronavirus Disease (COVID-19) Dashboard*. <https://covid19.who.int/>
- [3] CBC, retrieved January 21, 2021. *Here's the COVID-19 vaccine rollout plan, province by province*. <https://www.cbc.ca/news/canada/covid19-vaccine-rollout-plans-canada-1.5836262>
- [4] Government of Canada, retrieved January 21, 2021. *Vaccines and treatments for COVID-19: Vaccine rollout*. <https://www.canada.ca/en/public-health/services/diseases/2019-novel-coronavirus-infection/prevention-risks/covid-19-vaccine-treatment/vaccine-rollout.html>
- [5] Pelly, Lauren, retrieved January 21, 2021. *Health-care workers lining up for COVID-19 vaccine, but some warn of 'real troubles' with hesitancy*.
<https://www.cbc.ca/news/canada/toronto/health-care-vaccines-covid-19-hesitancy-1.5852119>