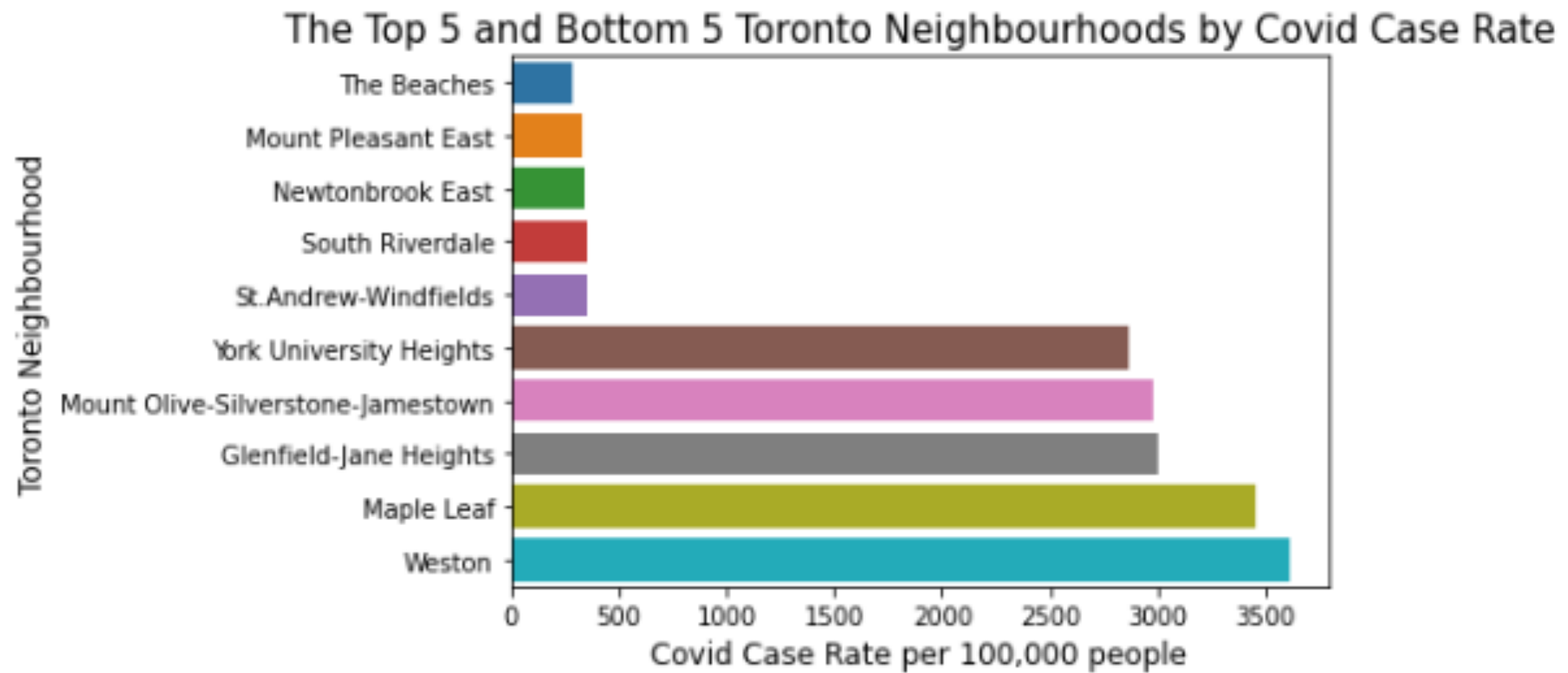# Predicting Covid rates in Toronto Neighborhoods using Linear Regression
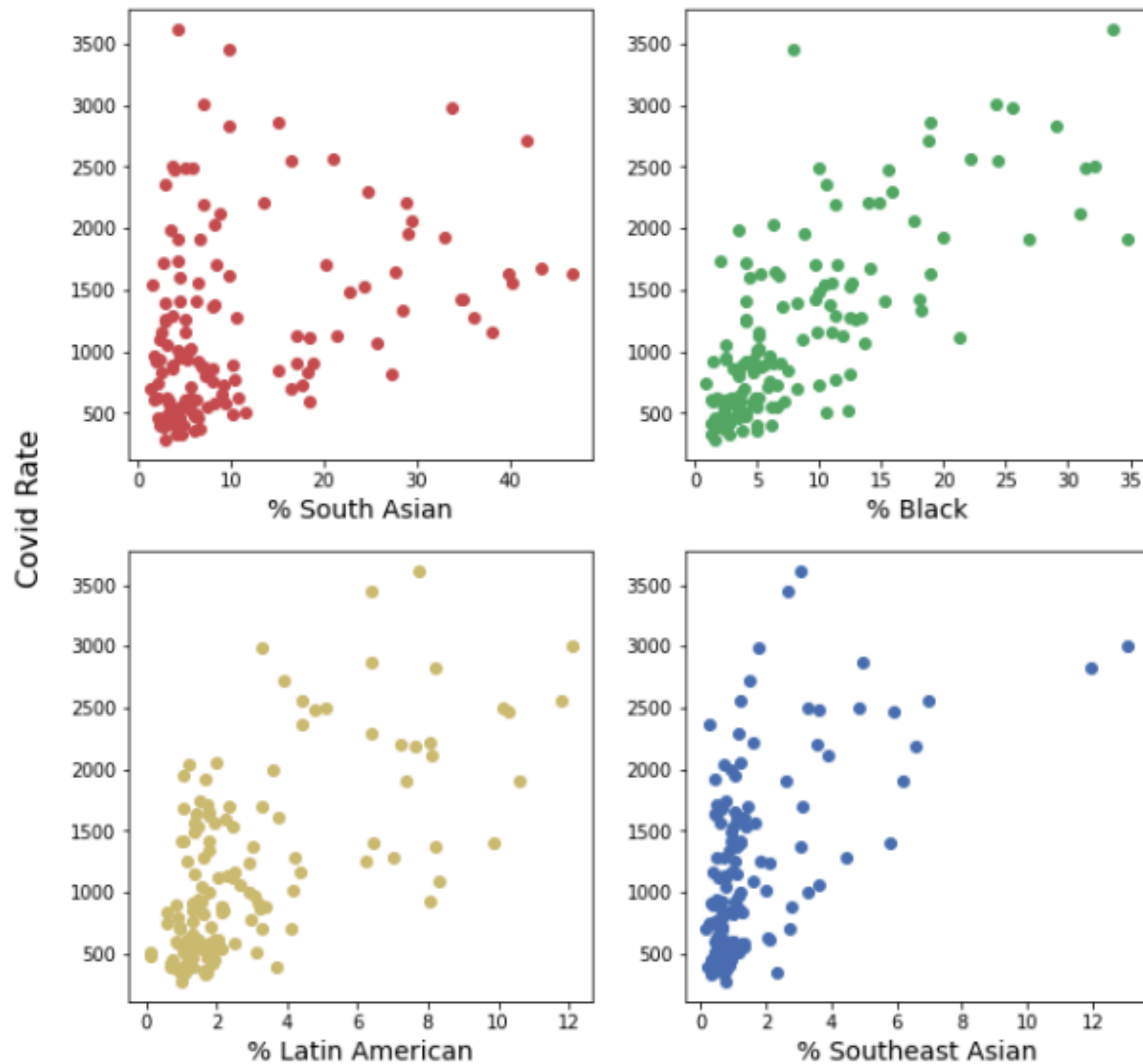
Alexei Marcilio
GBC

- ❖ Data from Toronto's open data portal

- ❖ Two files combined, covid rates and Census Data

- ❖ Over 2,300 potential features

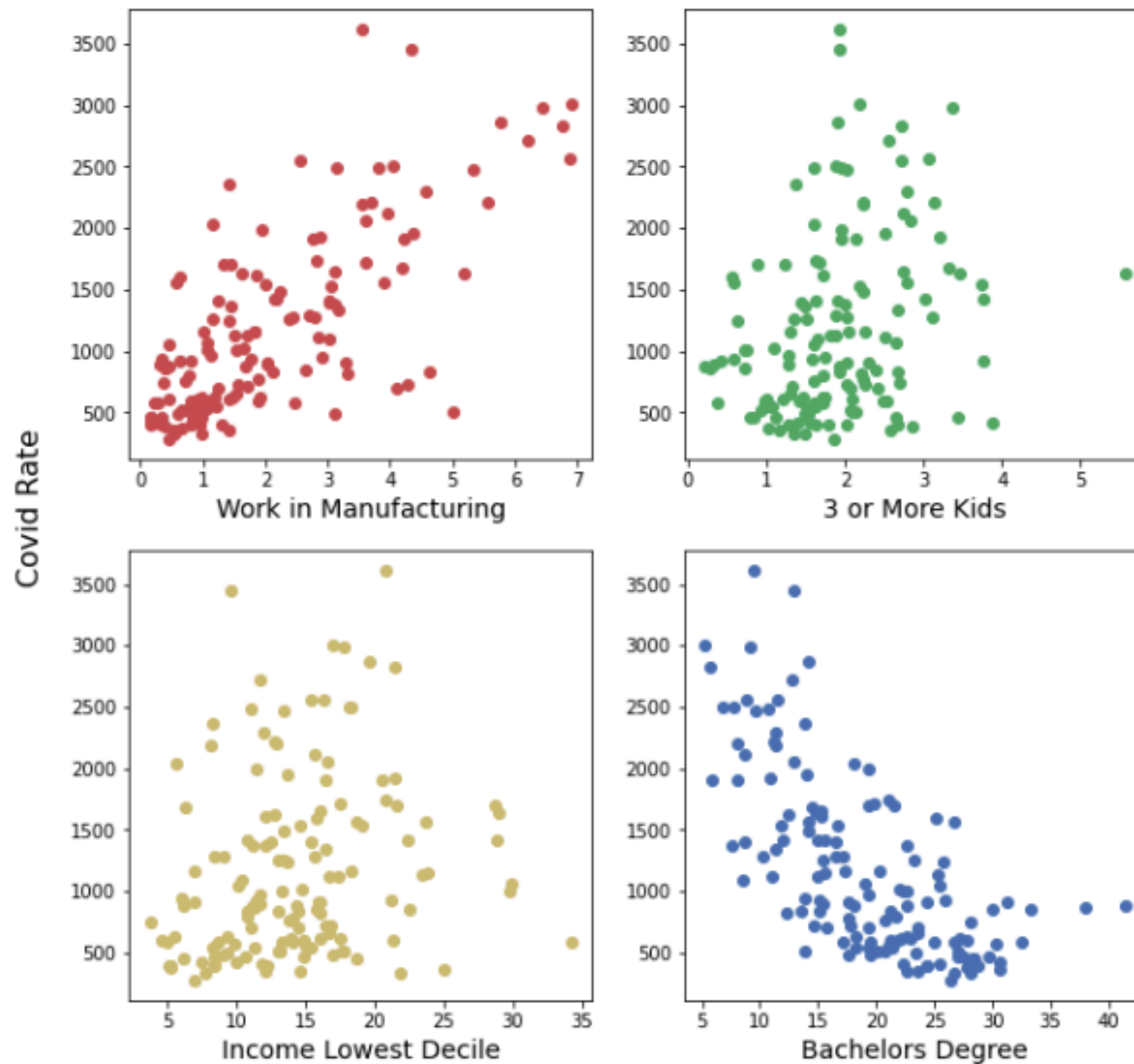The Top 5 and Bottom 5 Toronto Neighbourhoods by Covid Case Rate

Rates Vary

Covid Rates of Toronto Neighborhoods vs Percent of certain Races

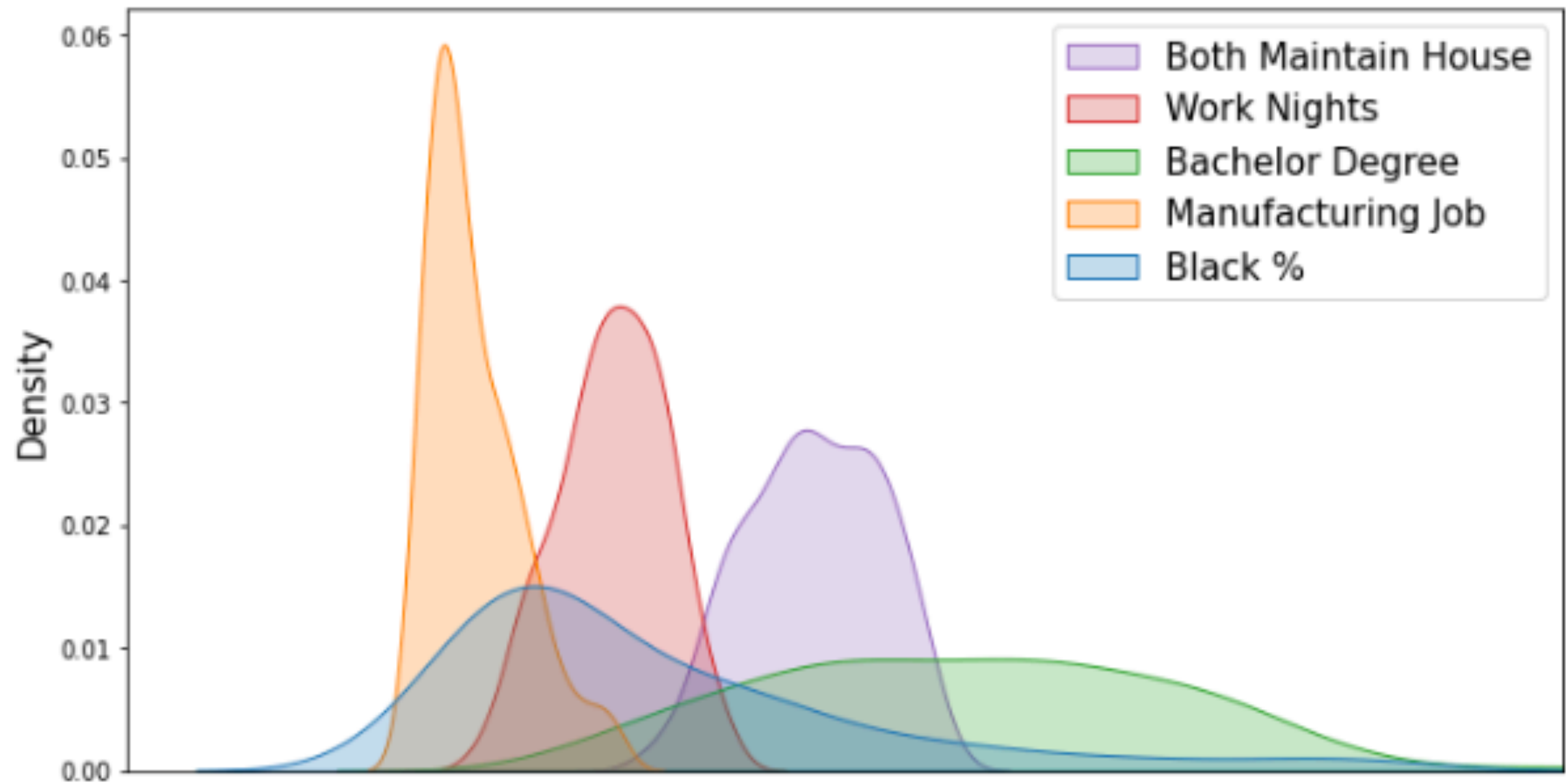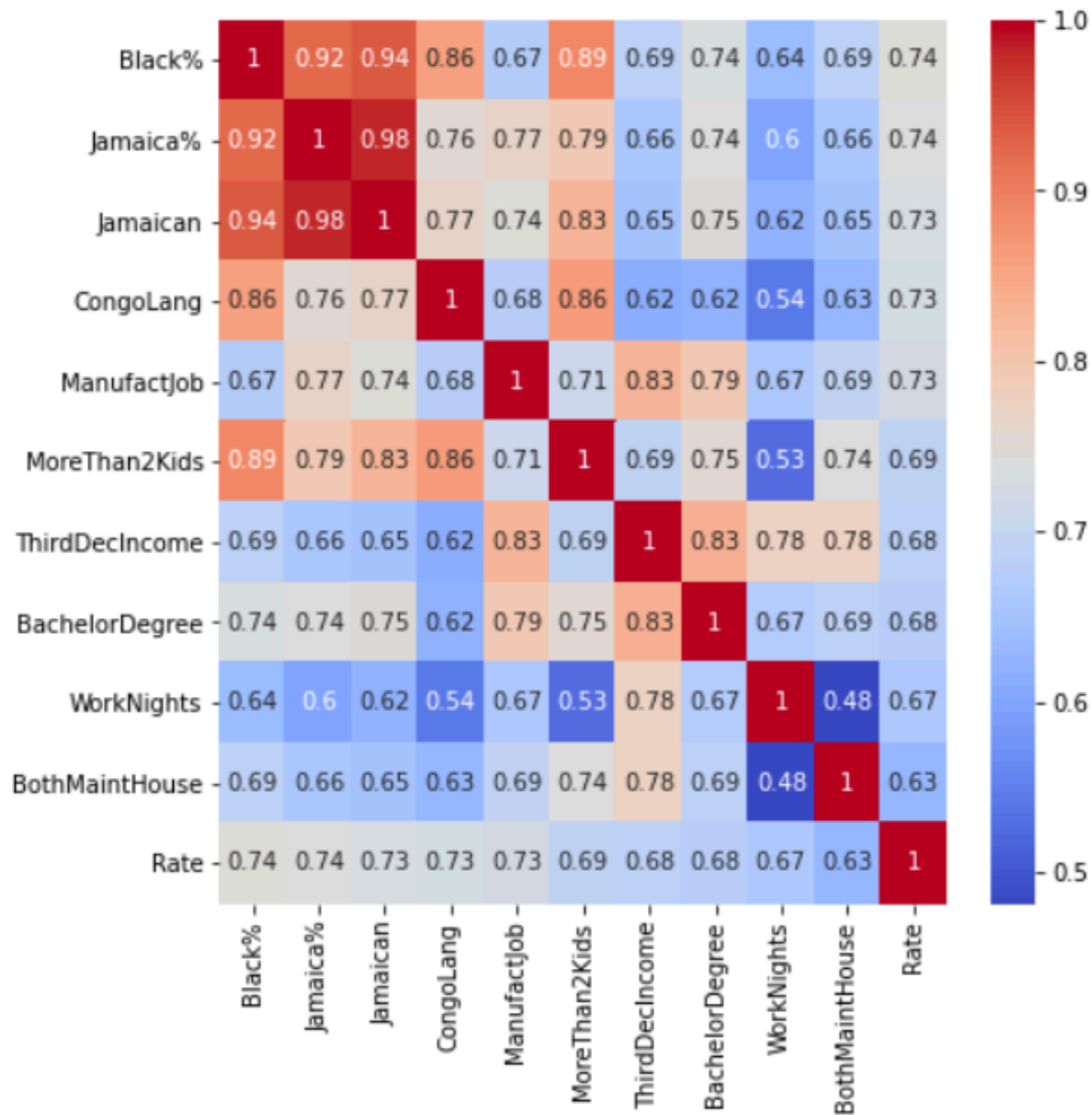Covid Rates of Toronto Neighborhoods vs Percent of Demographic Factors

| index | Rate per 100,000 people | Category | Characteristic |
|---|---|---|---|
| Col_1269 | 0.744447 | Visible minority | Black |
| Col_1105 | 0.741047 | Immigration and citizenship | Jamaica |
| Col_1377 | 0.731488 | Ethnic origin | Jamaican |
| Col_329 | 0.727007 | Language | Niger-Congo languages |
| Col_1855 | 0.726108 | Labour | 9 Occupations in manufacturing and utilities |
| Col_105 | 0.690762 | Families, households and marital status | 3 or more children |
| Col_1049 | 0.684972 | Income | In the third decile |
| Col_1635 | 0.677412 | Education | Bachelor's degree |
| Col_1907 | 0.666669 | Journey to work | Between 12 p.m. and 4:59 a.m. |
| Col_1594 | 0.629968 | Housing | 2 household maintainers |

# Choosing factors

Distributions of Our Predictors

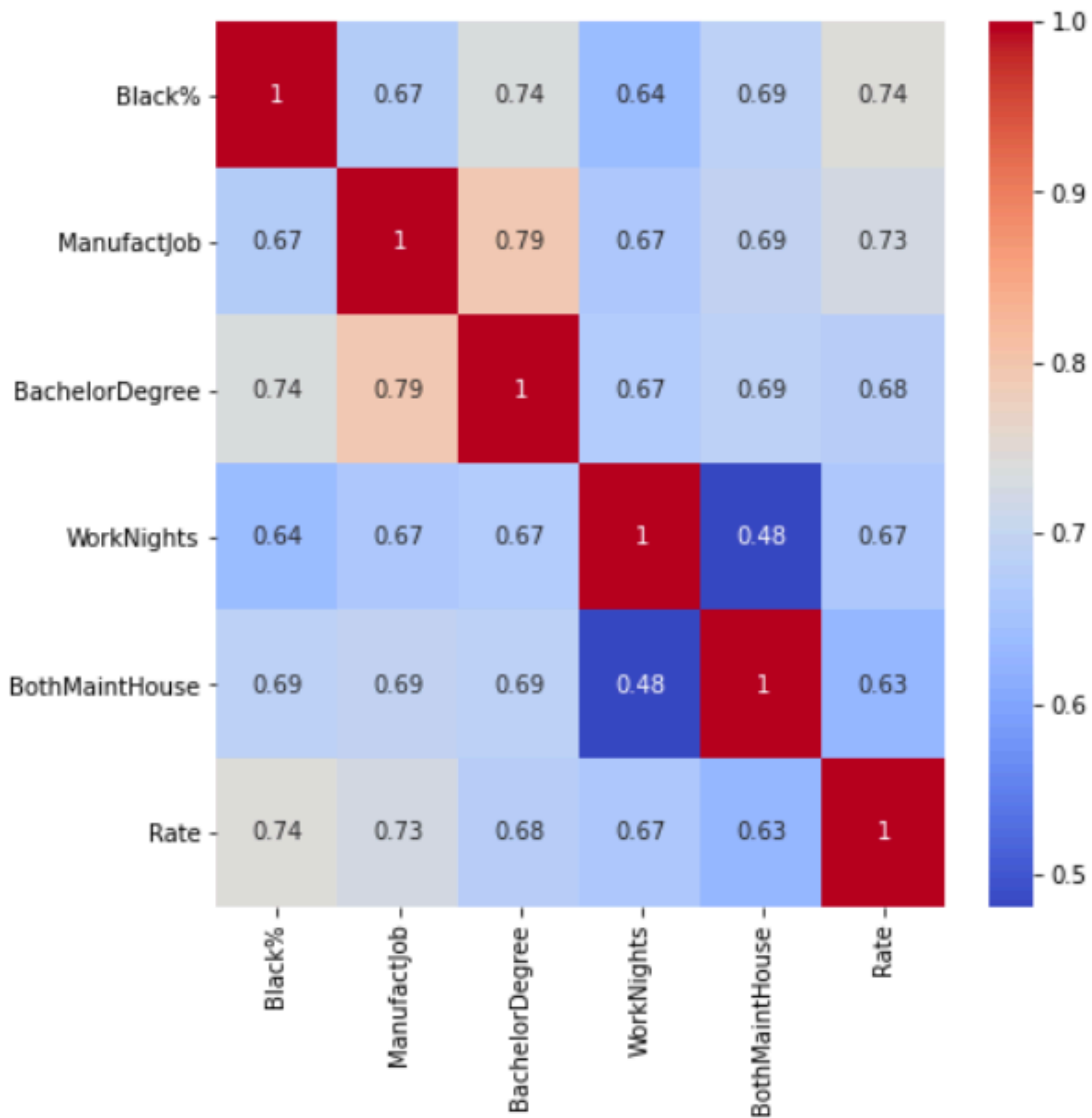Actual vs. Predicted results

Let's look at the results of the model. The $R^2$ **value is 0.43**. It's lower than simply using one predictor.

```
The R-squared value is:        43.47
The Root MSE is:               330.03675517892543
The Intercept is:              1212.2741827886584
```

Boxplots of the Features showing Outliers

```
The R-squared value is:      56.82
The Root MSE is:             432.34644674183556
The Intercept is:            1031.6761147565676
```

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
| Black% | 34.9662 | 11.486 | 3.044 | 0.003 | 12.228 | 57.704 |
| ManufactJob | 132.2317 | 38.480 | 3.436 | 0.001 | 56.057 | 208.406 |
| BachelorDegree | 7.6066 | 8.356 | 0.910 | **0.364** | -8.934 | 24.147 |
| WorkNights | 95.0617 | 28.633 | 3.320 | 0.001 | 38.379 | 151.744 |
| BothMaintHouse | -16.7414 | 16.746 | -1.000 | **0.319** | -49.892 | 16.409 |

The p-values indicate that there are potentially two features that do not contribute to the model. Let's remove one at a time and check the results. Here's the values after one predictor is removed.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
| Black% | 33.9600 | 11.425 | 2.972 | 0.004 | 11.345 | 56.575 |
| ManufactJob | 122.1344 | 36.821 | 3.317 | 0.001 | 49.249 | 195.019 |
| WorkNights | 94.1754 | 28.597 | 3.293 | 0.001 | 37.570 | 150.781 |
| BothMaintHouse | -3.8738 | 8.974 | **-0.432** | 0.667 | -21.637 | 13.890 |

| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
| Black% | 36.0324 | 10.333 | 3.487 | 0.001 | 15.580 | 56.484 |
| ManufactJob | 126.8286 | 35.063 | 3.617 | 0.000 | 57.429 | 196.228 |
| WorkNights | 83.0461 | 12.330 | 6.735 | 0.000 | 58.641 | 107.452 |

Now all our features are significant. Let's check the $R^2$ value and the MSE of the new model.

```
The R-squared value is:        60.52
The Root MSE is:               413.42407505558594
The Intercept is:              1038.1824194302878
```

- ❖ Linear Model is a good choice

- ❖ More study would be interesting

- ❖ Rate = 1038 + 157.7 * Black% + 140.8 * ManufactJob + 170.1 * WorkNights

- ❖ Racial Differences - underlying health, dense neighborhoods, lower % can work at home

Percent of those with Bachelor Degrees in Toronto Neighborhoods vs. Covid Rate per 100,000