

Examining the relationship between socioeconomic variables and air pollution in Toronto *

Farah Chin 400229991
Talat Hakim 400315290
Nathan Nadeau 400342430
Cindia Dao-Vu 400319161
Ifra Awan 400261667
Oliver Lawless 400271127

abstract goes here; some extra spiel added here

Keywords: pollution, demographics, spatial analysis, transportation

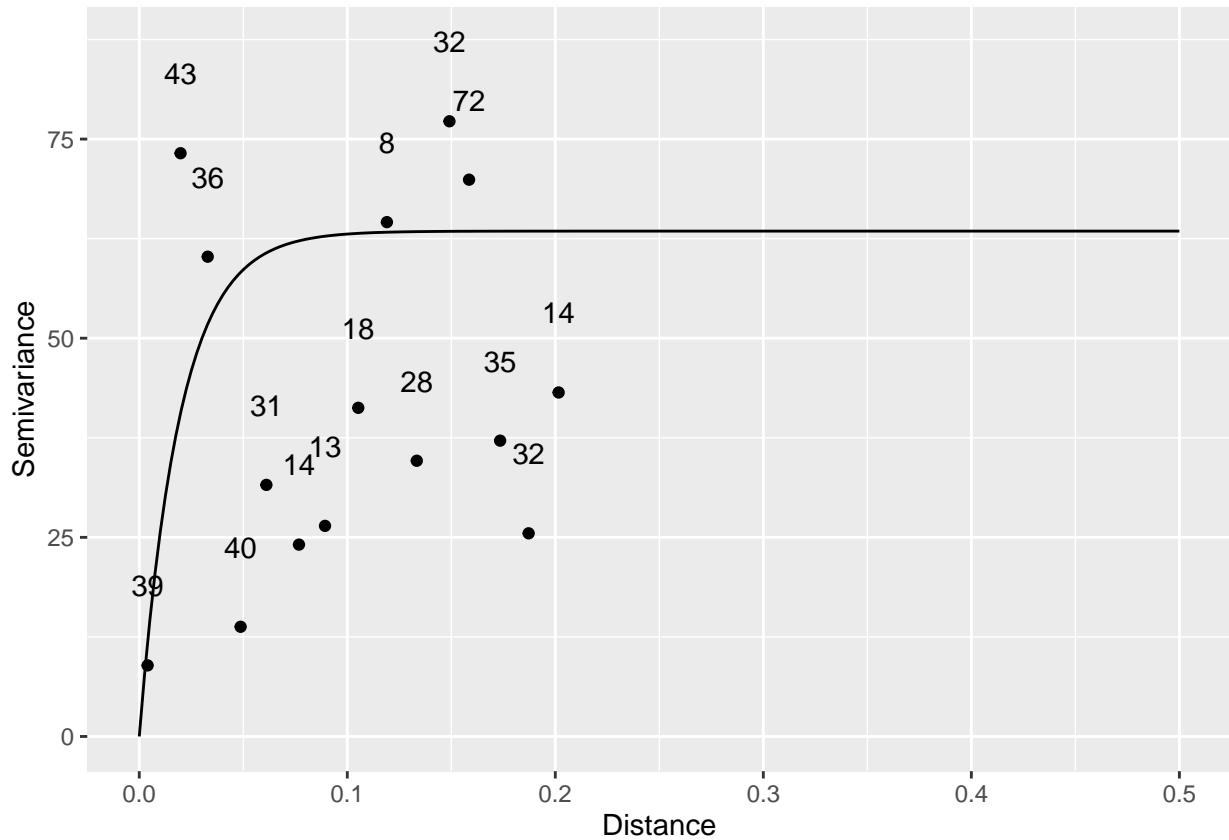
Introduction

This study analyzes the correlations between emission-based air pollution and socioeconomic spatial phenomena within the City of Toronto. It examines the predominant air quality factors responsible for the emission of fine particulate matter (PM2.5). The aim of this study is to address variations in air pollution statistics throughout the city. Generally, air quality tends to be worse in areas of lower socioeconomic status (SES) where people are forced to live near major sources of contamination such as major roads and factories. This paper contains a literature review to set a foundation for the analysis. Using air quality data from several sensors, we interpolated fine particulate matter concentrations across Toronto. Then, using data collected from the Canadian census we performed a linear regression analysis to determine the correlation between the air pollution of an area and the socioeconomic factors. Next, the statistics were summarized and interpreted to determine which socioeconomic factors are the most correlated with poor air quality. Finally, the paper concludes with some recommendations for future correlation analysis and steps the City of Toronto should take to manage air pollution.

*Paper submitted to complete the requirements of ENVSOCITY 4GA3 Applied Spatial Statistics; with additional edits by Antonio Paez for this version.

Background

Data



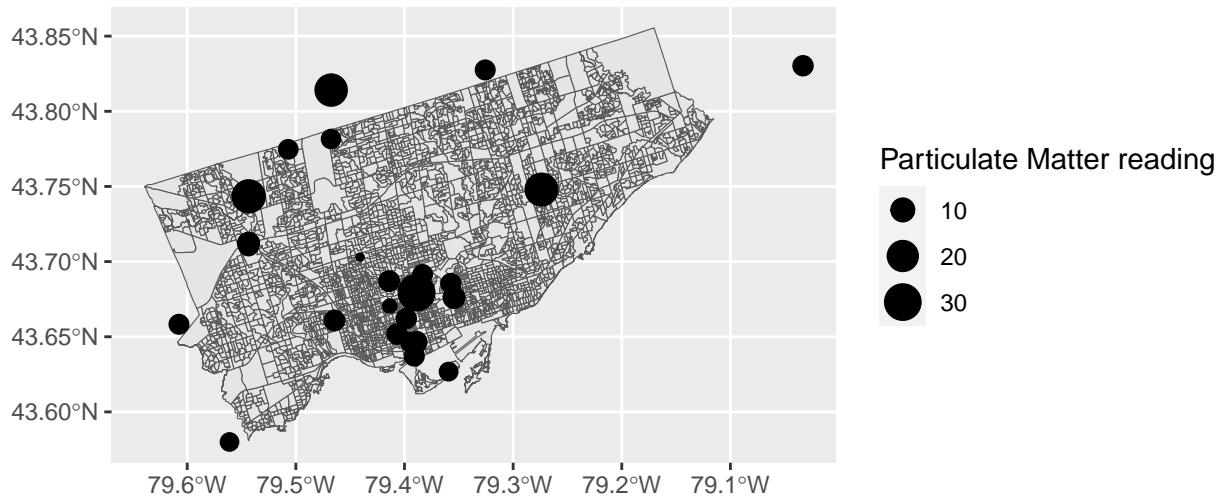
Toronto Census Data

Go to <https://mountainmath.github.io/cancensus/index.html> and follow instructions for setting up CensusMapping account and generating API key to retrieve data.

More info: <https://mountainmath.github.io/cancensus/articles/cancensus.html>

Set your cache to the data folder:

```
## [1] "C:/Users/Farah Chin/Documents/McMaster/ENVS0CTY 4GA3/Air-Pollution-Correlates-4GA3/data"
```

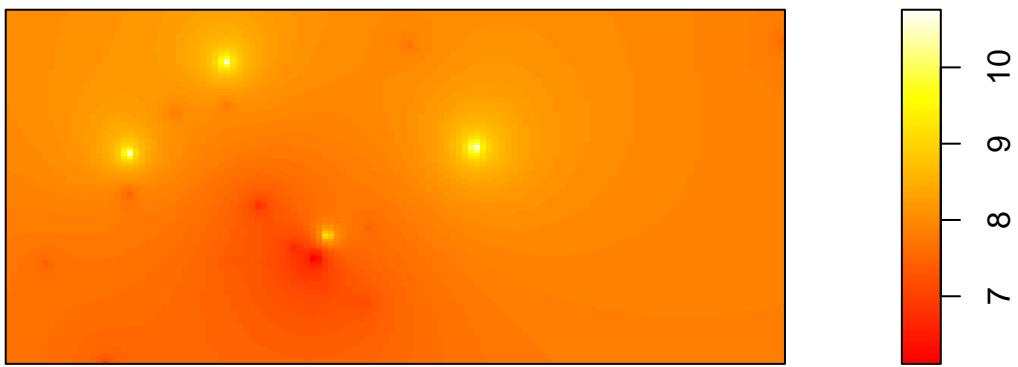


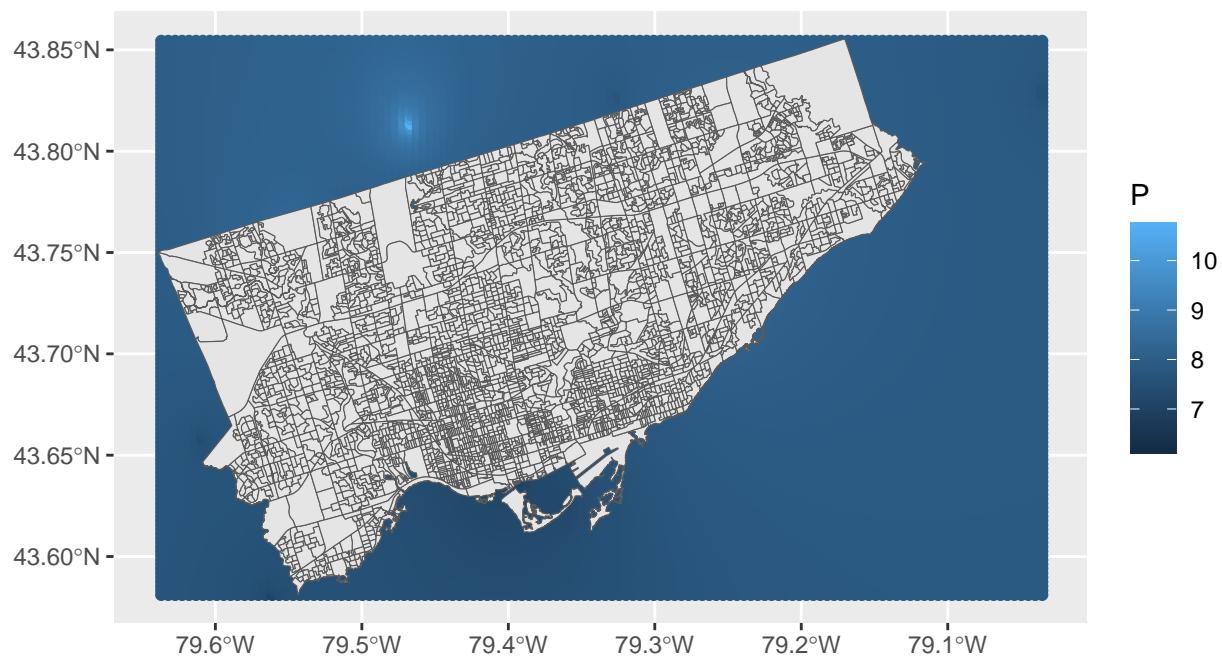
Here we use Leave One Out cross validation to determine the optimal power: <https://www.statology.org/leave-one-out-cross-validation/>

See also: https://rpubs.com/Dr_Gurpreet/interpolation_idw_R

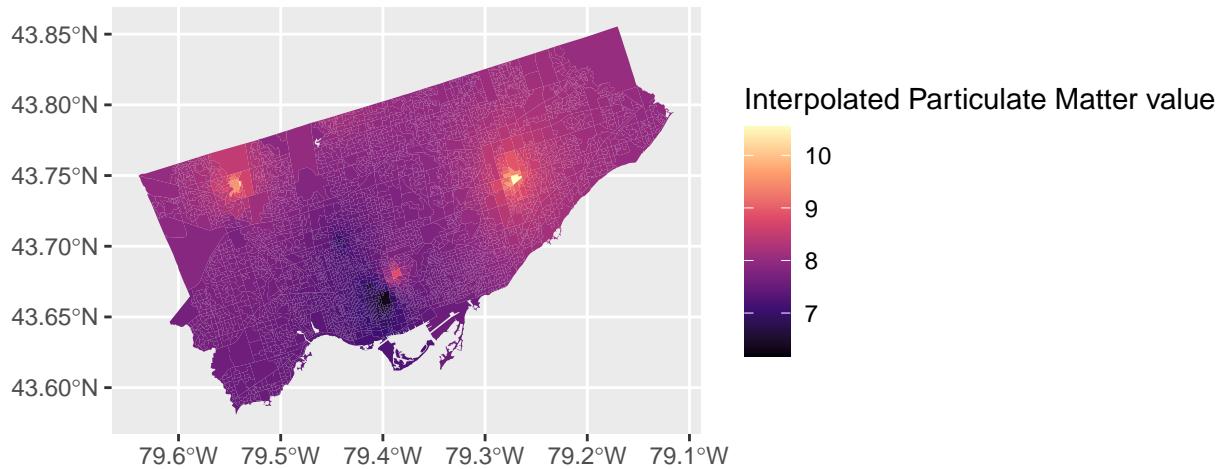
```
## [1] 0.471
```

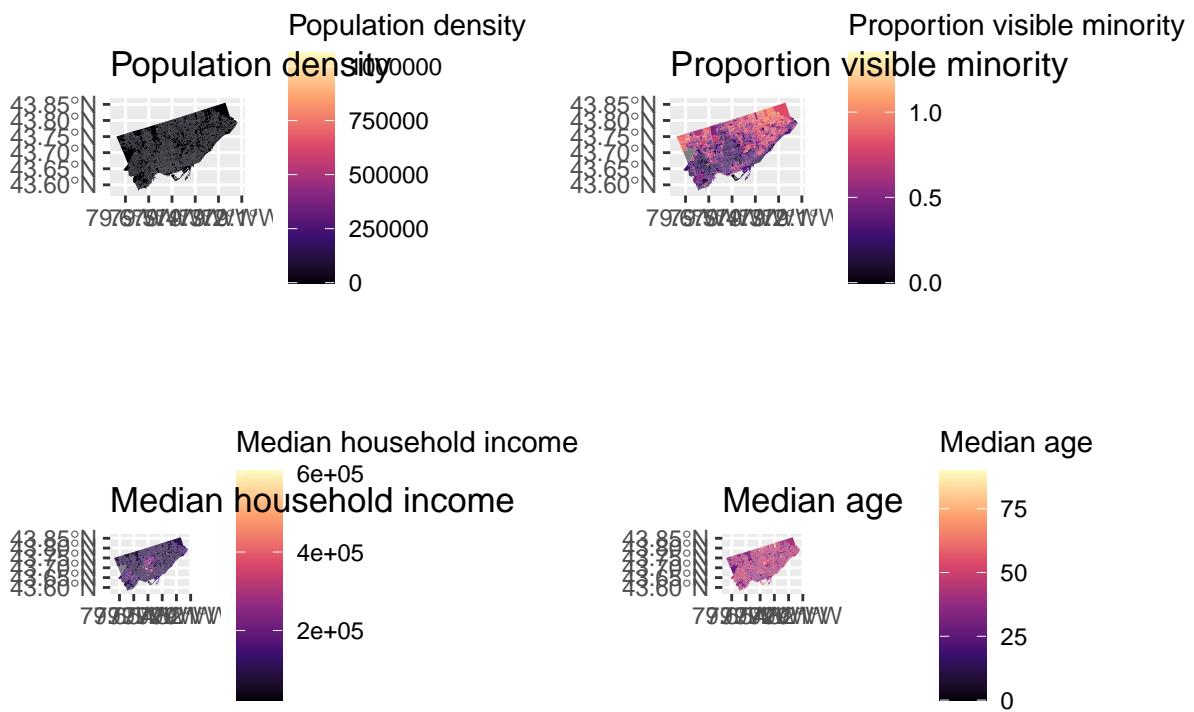
P.idw_pixels

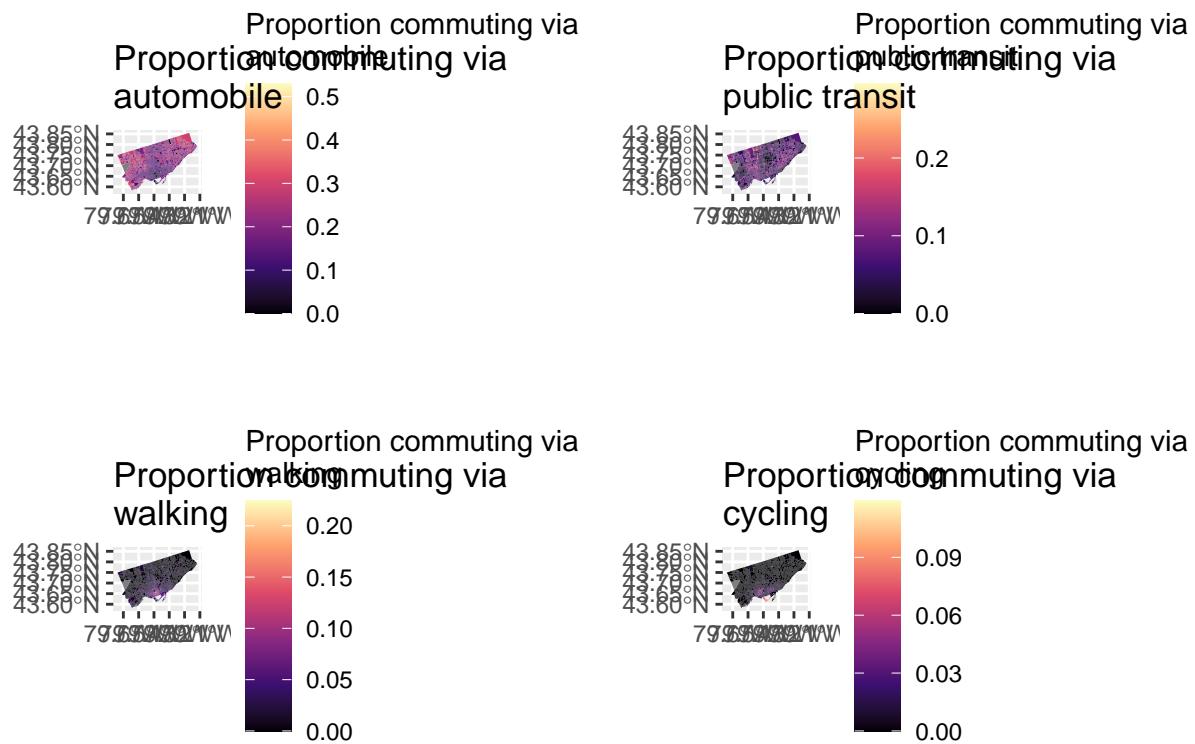


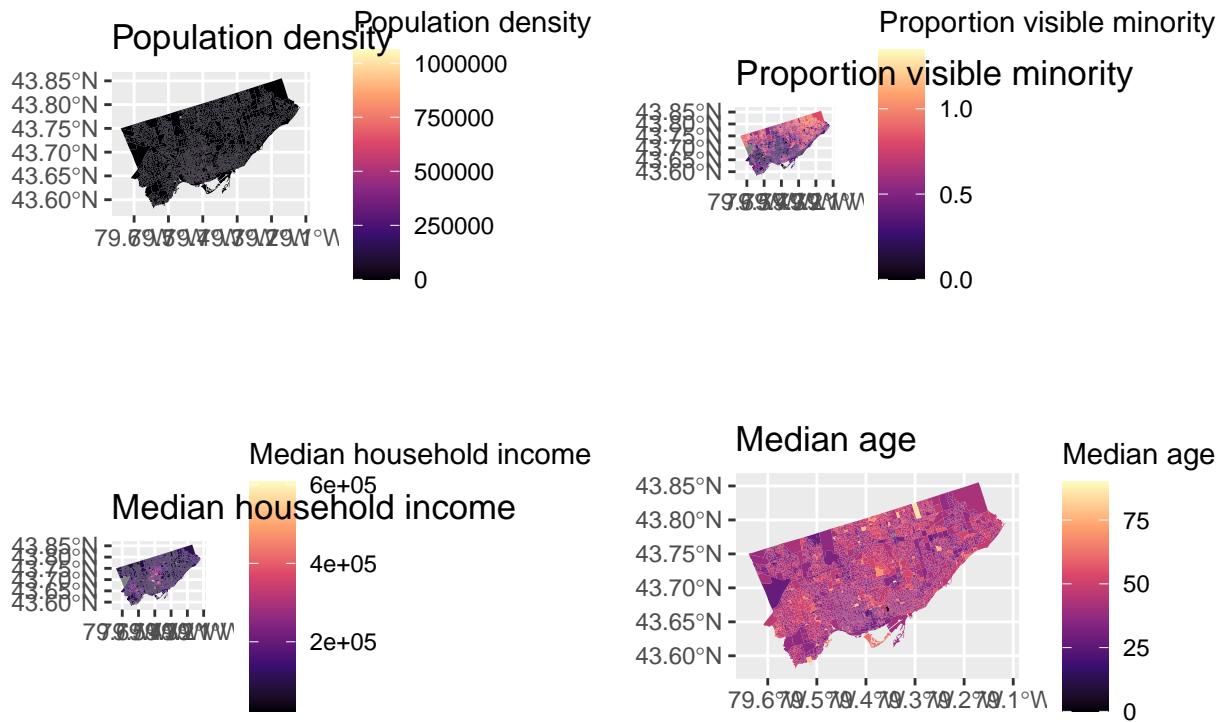


Interpolated Particulate Matter Readings





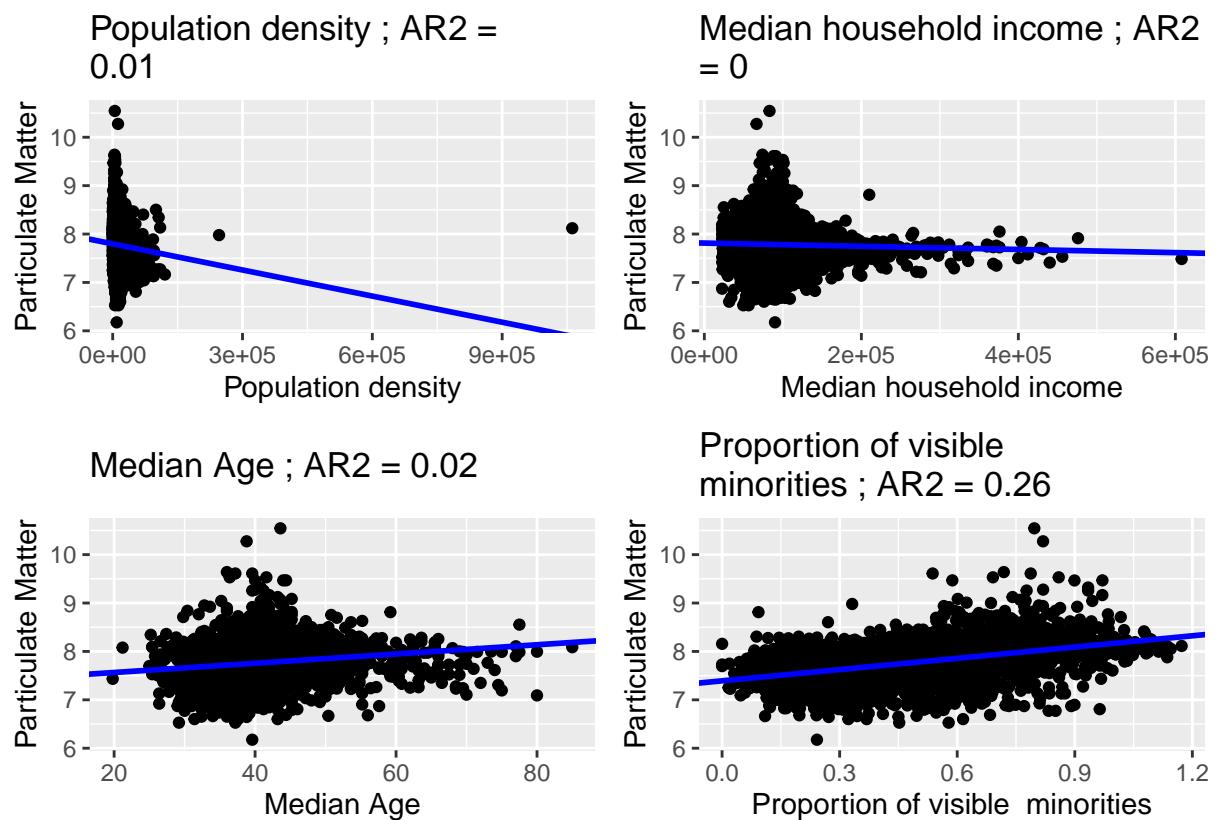


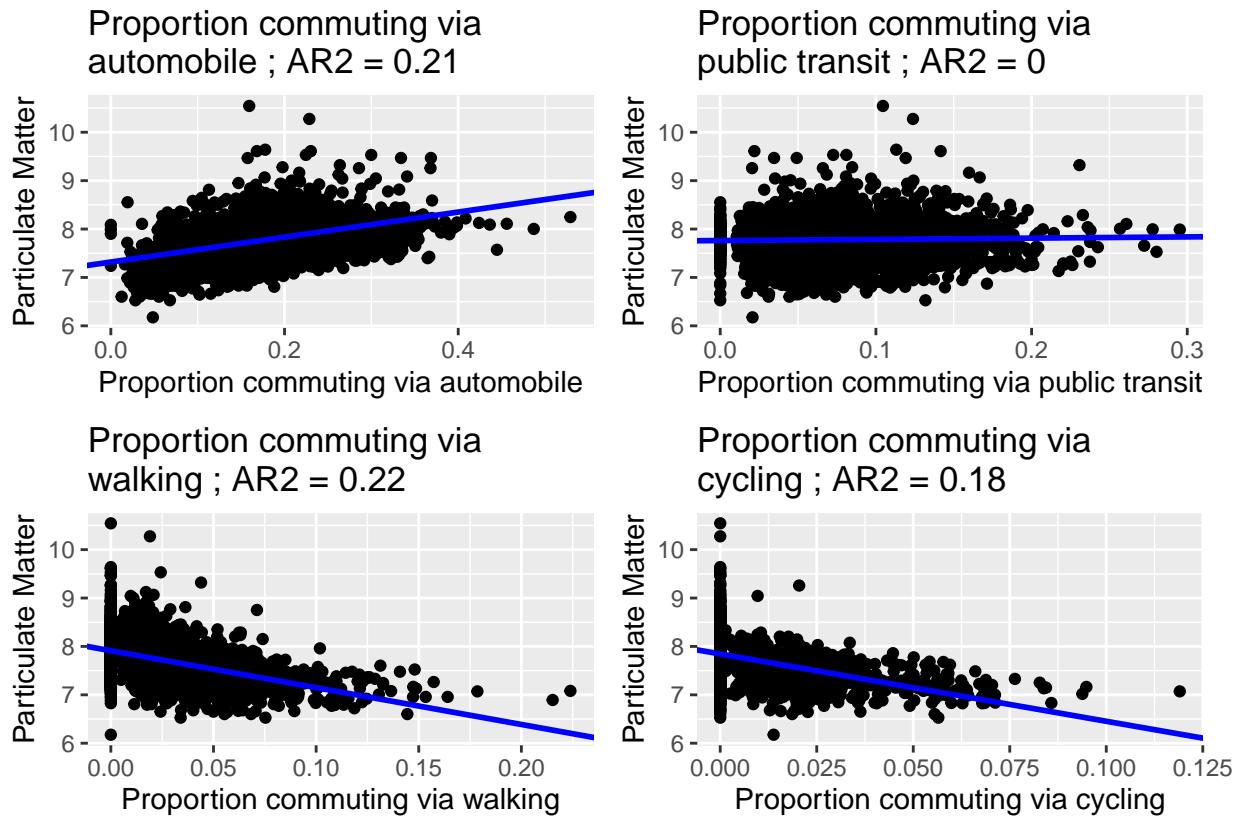


```

##      Variable
## [1,] "Population density"
## [2,] "Median household income"
## [3,] "Median age"
## [4,] "Proportion visible minority"
## [5,] "Proportion commuting via automobile"
## [6,] "Proportion commuting via public transportation"
## [7,] "Proportion commuting via walking"
## [8,] "Proportion commuting via cycling"
##      p-value
## [1,] "2.56077593971829e-68"
## [2,] "0"
## [3,] "1.4150620141619e-226"
## [4,] "0"
## [5,] "0"
## [6,] "3.66842309312586e-273"
## [7,] "0"
## [8,] "0"

```





To further investigate the relative importance of each covariate, linear regression was performed using all eight covariates. It is also possible that this might result in a better fit than models involving only single covariates.

```
##
## Call:
## lm(formula = P ~ population_density + minority_proportion + median_age +
##     median_hh_income + automobile_prop + public_tr_prop + walking_prop +
##     cycling_prop, data = toDAs.NAomit)
##
## Residuals:
##      Min        1Q        Median        3Q       Max
## -1.33640 -0.17765 -0.01834  0.16016  2.50336
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)             7.051e+00  4.948e-02 142.520 < 2e-16 ***
## population_density -4.722e-07  2.383e-07 -1.981  0.04762 *
## minority_proportion  6.684e-01  2.225e-02 30.039 < 2e-16 ***
## median_age            7.913e-03  8.249e-04  9.593 < 2e-16 ***
## median_hh_income    3.541e-07  1.250e-07  2.832  0.00465 **
## automobile_prop      7.975e-01  8.085e-02  9.863 < 2e-16 ***
## public_tr_prop      -4.072e-01  1.287e-01 -3.164  0.00157 **
## walking_prop         -3.683e+00  2.314e-01 -15.914 < 2e-16 ***
## cycling_prop         -5.215e+00  4.412e-01 -11.819 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

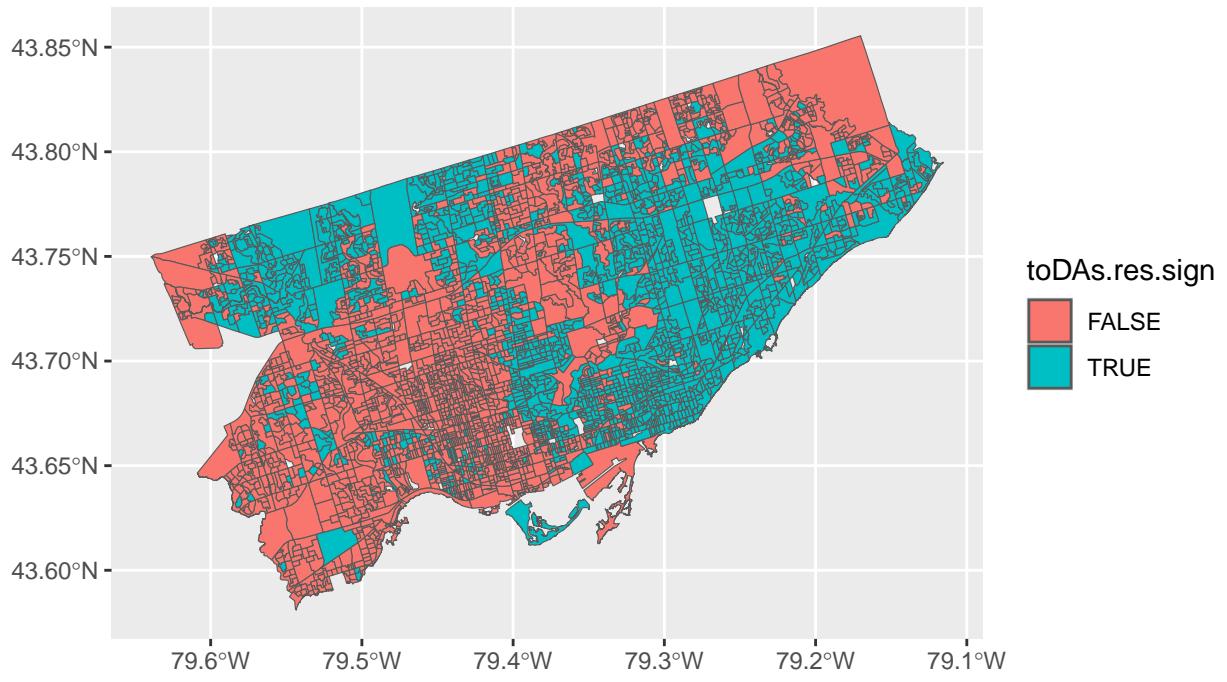
```

##
## Residual standard error: 0.29 on 3666 degrees of freedom
## Multiple R-squared:  0.4841, Adjusted R-squared:  0.483
## F-statistic:  430 on 8 and 3666 DF,  p-value: < 2.2e-16

##
## Call:
## lm(formula = P ~ minority_proportion + median_age + automobile_prop +
##      walking_prop + cycling_prop, data = toDAs.NAomit)
##
## Residuals:
##    Min      1Q  Median      3Q     Max
## -1.32615 -0.17306 -0.01627  0.16371  2.49662
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 7.0361088  0.0415344 169.40  <2e-16 ***
## minority_proportion 0.6154436  0.0200923  30.63  <2e-16 ***
## median_age       0.0088286  0.0007926  11.14  <2e-16 ***
## automobile_prop   0.8531586  0.0800063  10.66  <2e-16 ***
## walking_prop      -3.8315499  0.2294299 -16.70  <2e-16 ***
## cycling_prop      -5.3605210  0.4406873 -12.16  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2912 on 3669 degrees of freedom
## Multiple R-squared:  0.4795, Adjusted R-squared:  0.4788
## F-statistic: 675.9 on 5 and 3669 DF,  p-value: < 2.2e-16

```

Despite using only the covariates with the lowest p-value, the Adjusted R-squared value decreases slightly and the residual standard error increases. This indicates that reducing the number of variables in the model does not in fact improve the fit. Therefore, for general prediction it might be best to use all 8 of the selected variables.



Words

Methods

Describing some methods.

Results and discussion

Here we will discuss the results of the analysis.

References