# An Analysis on Commercial Real Estate Property Assessment and Property Tax in B.C.

Gian Carlo Di-Luvi   Mallory Flynn   Sophia Li   Vittorio Romaniello
STAT 550 Project Report
April 2020

Department of Statistics
University of British Columbia

# Summary

Property taxes are the single greatest operating expense for property owners in British Columbia, and mainly depend on two factors: property assessment—published each year in January—and municipal mill rates—published in April. Accurately projecting annual mill rates between January and April and future years' assessment values between May and December would allow businesses and individuals to budget for these expenses. For this purpose, multiple competing models are fitted based on a data set containing assessment values, mill rates, and other property information from 2016 to 2020 of over 200,000 properties. For predicting mill rates early in the year, a random forest regression model is found to be the most accurate. Likewise, another random forest regression model is also the most accurate for predicting future years' assessment values. Both models are incorporated in an interactive Shiny app that provides a straightforward user interface with a property-specific tax assessment.

# 1    Introduction

For property owners in British Columbia, property taxes are the single greatest operating expense. Mill rates are determined by individual municipalities and depend on a number of factors, including assessment values and municipal budgets. While property assessments are published each year in January, municipal mill rates are not determined until April of the same calendar year, making it difficult for property owners to budget for this large expense. Thus, accurately projecting the annual mill rates between the months of January and April and the next year's assessment values between May and December is a value-add to real estate consulting services. Furthermore, it also assists property owners in financial planning and fund allocation. The objective of this written report is to build accurate statistical predictive models for both future municipal mill rates and property-specific assessment values. This will be done for tax class codes 1 (residential), 5 (light industrial), and 6 (commercial). Guided by the exploratory data analysis in Section 2.2, several modeling techniques are presented in Section 2.3. Mill rate and assessment predictions using these models are available through a straightforward user interface via a Shiny app, which is discussed in the same section. Finally, Section 3 includes some concluding remarks.

# 2    Statistical analysis

## 2.1    Data

Han et al. (2020) use the same dataset and offer a thorough description of the data in their analysis. In this report, however, all B.C. municipalities are included in analysis. The variables are summarised in Table 1. Missing mill rate values were imputed by using the average municipality mill rate for the corresponding year and tax class code. This is sensible because mill rates are constant across municipalities, and so a more complex imputation technique is not necessary.

An important feature of the data is that spatial relationships exist within municipalities for assessment values and between municipalities for mill rates. In other words, it is not sensible to assume independence in neither assessment values nor mill rates between different properties because geographical location (and thus proximity) affects these variables. This significantly influences possible model choices as the independence assumption is required for several standard modeling techniques. Furthermore, the data set lacks appropriate spatial information that would allow for geographical variables to be taken into account.

| Variable | Type |
|----------|------|
| Mill rate | Continuous |
| Total assessment | Continuous |
| Total land assessment | Continuous |
| Total improvement assessment | Continuous |
| Tax class code | Categorical. One of 1 (residential), 5 (industrial), 6 (commercial) |
| Municipality | Categorical |
| Year | Discrete. Values of 2016, 2017, 2018, and 2019 |

Table 1: Brief description of the variables included in the data set.

## 2.2  Exploratory data analysis

The exploratory data analysis (EDA) further investigates aspects of the EDA done by Han et al. (2020). Both assessment values and mill rates have right-skewed distributions, and so they were log-transformed for visualization purposes. A visual inspection of Figure 1a shows that mill rates differ drastically from tax class to tax class, with industrial properties having the highest mill rates. However, neither mill rates nor assessment values seem to vary much over time, which suggests that current values of these variables will be a good indicator of their (short-term) future values. Furthermore, there does not appear to be a clear relationship between assessment values and mill rates when tax class is taken into account—except in the residential tax class, where a slightly negative relationship is present between the logarithms of average municipal assessment values and municipal mill rates. This remains true when only taking into account values from 2020, as can be seen in Figure 1b.



(a) Scatterplot with year in different colors.    (b) Scatterplot only for 2020.
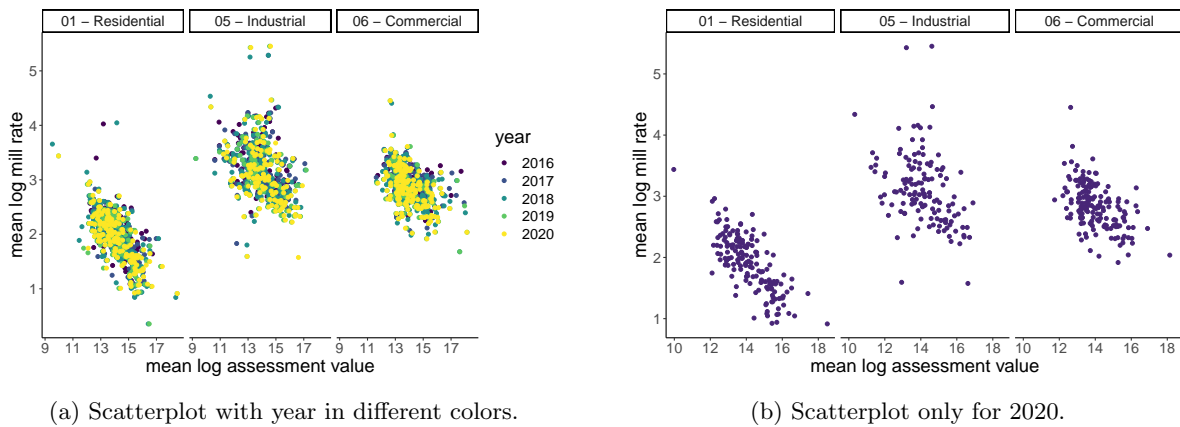
Figure 1: Average log assessment values against log mill rates by municipality across the three different tax classes of interest.

The fact that mill rates remain relatively constant throughout the years is further verified in Figure 2a, where industrial properties are again seen to have the largest mill rates on average. However, Figure 2b suggests that assessment values do not follow this pattern. Municipality plays a role in a property's assessments, something taken into account in the modeling process.

Finally, Figure 3 shows that the distributions of both mill rates and assessment values are indeed right-skewed (the values shown are in logarithmic scale), and further confirms the stark difference in mill rates between the three tax classes. Notably, Figure 3b suggests that this difference is not present in assessment

values, which indicates that tax class may not be that relevant for predicting future years' assessment values.
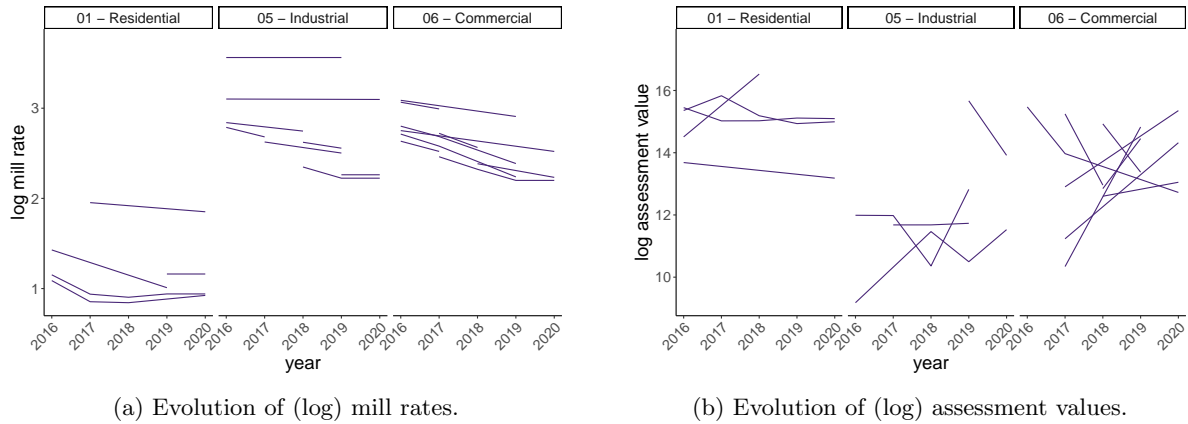


(a) Evolution of (log) mill rates.



(b) Evolution of (log) assessment values.

Figure 2: Evolution of target variables over the years for some randomly-selected properties.



(a) Distribution of (log) mill rates.
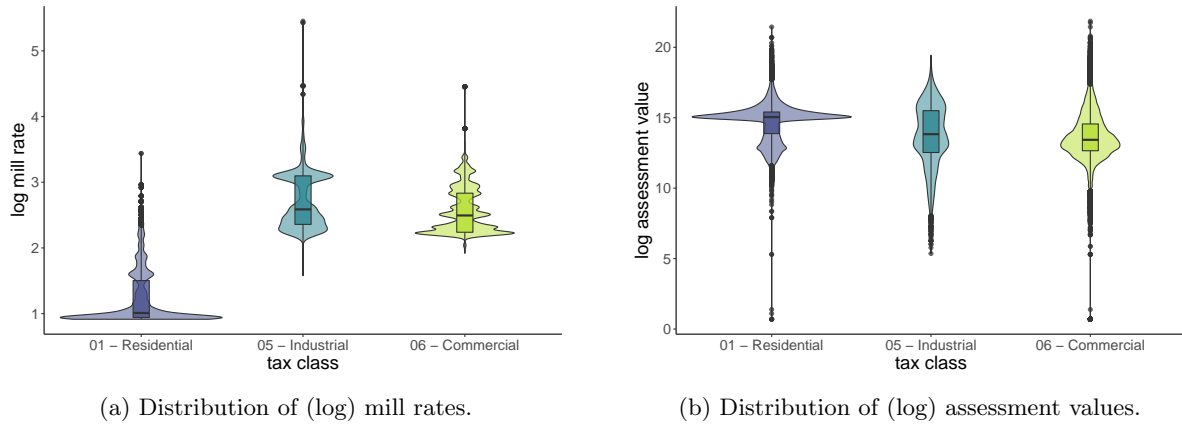


(b) Distribution of (log) assessment values.

Figure 3: Violin and boxplots of (log) mill rates and (log) assessment values of individual properties across the three tax classes of interest, taking into account only the year 2020.
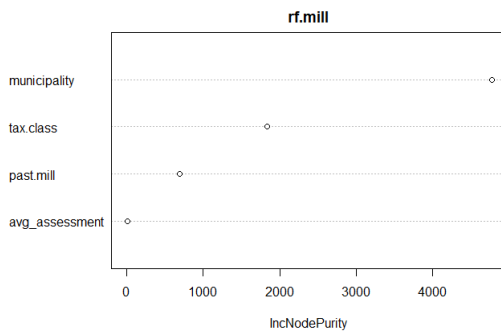
## 2.3 Modeling

This section describes the models used to predict mill rate and assessment values. Each of the following models were trained on a subset of 75% of the data and tested on the remaining 25%.

### 2.3.1 Random forest regression

The random forest (RF) algorithm is an ensemble learning method for classification and regression that operates by constructing numerous decision trees and choosing the majority class output (classification) or the average prediction (regression) given by the individual decision trees. Decision trees are trained on different samples (with replacement) of the data, and only a random subset of the covariates is used to define the split at each node, based on a recursive partitioning strategy. RF is preferred over a single decision tree as it overcomes the problem of over-fitting. The following models were constructed using the 'randomForest' package in R (Liaw & Wiener, 2002).

A RF model was fitted to predict mill rate using tax class, municipality, past years' mill rate, and current years' total assessment as independent variables. Another RF model was similarly constructed using next year's assessment value as the dependent variable, and municipality, current years' mill rate, tax class, and current years' total assessment as covariates. The number of trees used for both models was 500, and all 4 predictors were evaluated at each node for splitting. The percent of total variability explained is 96.38 % for the mill rates model and 97.48 % for the assessment values model.

The importance of model covariates is indicated by a higher value of increase node purity (mean decrease accuracy). Municipality was deemed the most important variable in predicting mill rate (see Figure 4a), followed by tax class and past year's mill rate. However, as shown in Figure 4b, current years' assessment value is the most important variable in predicting next years' assessment value, followed by municipality and mill rate.



(a) Variable importance plot for mill rates.



(b) Variable importance plot for assessment values.

Figure 4: Variable importance plot based on increase in node purity for mill rate and total assessment value predictions.

Figure 5 shows the evaluation of the model's performance on a test dataset. The RF model tends to underestimate mill rates, probably due to the greater difficulty of the task. However, RF shows good performance in predicting next year's assessment values.



(a) Performance for mill rates prediction.



(b) Performance for assessment value prediction.

Figure 5: Models' performance on test set.

### 2.3.2  Mixed effects models and generalized estimating equation

Since the data set contains repeated measurements for both municipalities and properties over the years, it is reasonable to consider a linear mixed effect model (LME) for predicting mill rates and assessment values. Both random slope and random intercept for each municipality are needed based on Figure 2b, which depicts different initial values and rates of change for both assessment values and municipal mill rates. Based on the fixed effects shown in Figure 6b, the direction of change in assessment values fluctuates over time. For example, assessment values increased from 2016 to 2017 but decreased from 2018 to 2020. The median municipal assessment value is positively associated with current mill rate, as is past year's mill rate (Figure 6a). A negative relationship between mill rate and assessment values is observed in Figure 6b, indicating that lower mill rates are associated with higher assessment values.

```
Random effects:
 Groups        Name        Variance  Std.Dev. Corr
 municipality (Intercept) 7.023e-01 0.838028
              year        1.868e-06 0.001367 1.00
 Residual                 7.037e-01 0.838889
Number of obs: 141, groups:  municipality, 41

Fixed effects:
                      Estimate Std. Error t value
(Intercept)          1.451e+00  1.101e+00   1.318
as.factor(year)2018 -6.181e-01  2.080e-01  -2.972
as.factor(year)2019 -1.098e+00  2.303e-01  -4.766
as.factor(year)2020 -8.532e-01  2.596e-01  -3.287
tax.class05          1.728e+01  2.774e+00   6.228
tax.class06          1.125e+01  6.965e-01  16.152
avg_assessment       3.583e-07  3.668e-07   0.977
past.mill            2.979e-01  4.012e-02   7.425
```

(a) Summary output for predicting mill rates.

```
Random effects:
 Groups        Name        Variance  Std.Dev. Corr
 municipality (Intercept) 3.832e+15 61900409
              year        1.061e+09    32572 -1.00
 Residual                 3.825e+15 61846281
Number of obs: 753, groups:  municipality, 52

Fixed effects:
                      Estimate Std. Error t value
(Intercept)         -4.150e+06  9.212e+06  -0.450
as.factor(year)2017  7.074e+06  7.467e+06   0.947
as.factor(year)2018  1.229e+07  7.250e+06   1.695
as.factor(year)2019 -1.247e+07  7.304e+06  -1.707
as.factor(year)2020 -4.888e+06  1.152e+07  -0.424
tax.class05          1.006e+07  1.281e+07   0.786
tax.class06          1.308e+07  7.535e+06   1.735
total.assessment     1.028e+00  8.108e-03 126.731
mill.rate           -3.575e+05  4.714e+05  -0.758
```

(b) Summary output for predicting assessment values.

Figure 6: Models summary output.



(a) Diagnostics plots for LME model for mill rate.

(b) Diagnostics plots for LME model for assessment values.

Figure 7: Diagnostics plots.

5

Unfortunately, the diagnostics plots verifying model assumptions are subpar for both LME models, as can be seen in Figure 7. Residuals are not normally distributed, variance is heterogeneous, and random intercept is non-normal. Thus it does not appear that an analysis using this type of model is appropriate for this data.

### 2.3.3   Shiny app

A straightforward user interface for making predictions using the RF models for mill rate and assessment values was built by way of a shiny app, which can be found in the GitHub repository (Di-Luvi et al., 2020) or accessed on the web at https://malloryjflynn.shinyapps.io/STAT550_Real_Estate/.

The shiny app calculates predictions of the next years' mill rate or assessment value under the 'Estimate' tab of the main panel. This can be done using either a Property Identifier Code (PIC) or user inputs for municipality, tax class code, and previous year's mill rate (mill rate predictions) or current assessment value (assessment value predictions). To provide the user with a broader view of the trend in mill rate for the municipality selected, a plot of the mill rate for that municipality over the years (2016 - 2020) is visible under the 'Plot' tab of the main panel. When predicting assessment value using a PIC, a plot of the assessment value for that property over time will be displayed alternatively.

## 3   Conclusions

In this report we built models to predict mill rate and assessment value for properties in British Columbia's municipalities. Random forest was found to be the best fitting model, although it tended to overestimate mill rate predictions. In addition to this, we incorporated the predictive models in a shiny app that allows a user to perform predictions given a set of inputs and to visualize trends in mill rate and assessment value.

The main limitation of our analysis lies in the fact that the models assume independence of the data. However, as was argued in Section 2.1, mill rates and assessment values are highly dependent on spatial information. An attempt to obtain the properties' spatial information was made by transforming the addresses into spatial coordinates using an open source geocoding system provided by the British Columbia Government (see the documentation of Esmukov and Tigas (2018)). However, the open source geocoding system did not recognize some of the addresses in the dataset and provided incorrect coordinates for a large number of properties. Therefore it is recommended that (1) addresses be formatted according to the standards required by the B.C. government, or (2) a private geocoding system that may be able to find addresses using the current format in the data is used. Provided spatial information is available, several model alternatives (e.g. factor models) could be tested to improve prediction accuracy, which take spatial information into account.

# References

Di-Luvi, G. C., Flynn, M., Li, S., & Romaniello, V. (2020). *Stat450-550: Real estate consulting project.* Retrieved April 13, 2020, from https://github.com/STAT450-550/RealEstate

Esmukov, K., & Tigas, M. (2018). Geopy 2.0. https://github.com/geopy/geopy

Han, P., Lu, X., Liu, Y., & Wen, Y. (2020). *Stat 450 project: Real estate* (tech. rep.). Department of Statistics, University of British Columbia.

Liaw, A., & Wiener, M. (2002). Classification and regression by randomforest. *R News*, *2*(3), 18–22. https://CRAN.R-project.org/doc/Rnews/

# Appendix

## A   EDA code

```
1   # preamble ####
2   library(tidyverse)
3   library(readr)
4   library(readxl)
5   library(ggplot2)
6   ggplot2::theme_set(theme_classic())
7   library(viridis)
8
9   # data wrangling ####
10
11  # data import
12  columnnames <- readxl::read_xlsx("colnames.xlsx") %>%
13    colnames()
14
15  real.estate_full <- readr::read_csv("2016 - 2020 Raw.csv",
16                                      na = c("", "NA", "NULL", "NULL_1"),
17                                      col_names = columnnames,
18                                      col_types = "cicccccccccccccddddddiicccddddc")
19  # data wrangling
20
21  # assessments by PIC and year
22  assessments <- real.estate_full %>%
23    dplyr::select(PIC, Year, # relevant variables
24                  AssessedValueAmt, AssetTypeDesc) %>%
25    dplyr::rename(year = Year,
26                  assessment = AssessedValueAmt,
27                  assessment.type = AssetTypeDesc) %>%
28    dplyr::group_by(PIC, year, assessment.type) %>%
29    dplyr::summarise(assessment = sum(assessment)) %>%
30    dplyr::ungroup() %>%
31    tidyr::spread(assessment.type, assessment) %>%
32    dplyr::rename(improvement.assessment = Improvement,
33                  land.assessment = Land) %>%
34    dplyr::select(PIC, year, improvement.assessment, land.assessment) %>%
35    dplyr::mutate(total.assessment = improvement.assessment + land.assessment)
36
37
38  re <- real.estate_full %>%
39    dplyr::select(PIC, Year, AddressAssessorMunicipalityDesc, # relevant variables
40                  TaxClassCode, TaxOwingAmountTotalCalculated, TaxClassTaxRate) %>%
41    dplyr::rename(year = Year,
42                  municipality = AddressAssessorMunicipalityDesc, # human-readable names
43                  tax.class = TaxClassCode,
44                  tax = TaxOwingAmountTotalCalculated,
45                  mill.rate = TaxClassTaxRate) %>%
46    dplyr::filter(tax.class %in% c("01", "05", "06")) %>%  # relevant values for tax class
47    dplyr::distinct() %>%
48    dplyr::left_join(assessments, by = c("PIC" = "PIC", "year" = "year")) # add assessment
49
50
51  # data viz ####
52  # tax classes dictionary
53  tax.classes <- as_labeller(c(
54    `01` = "01 - Residential",
55    `05` = "05 - Industrial",
56    `06` = "06 - Commercial"
57  ))
58
59
60  # facet scatter plots with year
61  re %>%
62    dplyr::select(-PIC) %>%
63    dplyr::filter(!is.na(total.assessment), !is.na(mill.rate)) %>%
64    dplyr::group_by(year, municipality, tax.class) %>%
65    dplyr::summarise(total.assessment = mean(total.assessment), mill.rate = mean(mill.rate)) %>%
66    ggplot(aes(x = log(total.assessment), y = log(mill.rate), color = factor(year))) +
```

```r
 67    geom_point() +
 68    facet_wrap(.~tax.class, labeller = tax.classes) +
 69    labs(x = "log assessment value",
 70         y = "log mill rate",
 71         color = "year") +
 72    scale_color_viridis_d() +
 73    theme(text = element_text(size = 18))
 74  ggsave("RealEstate/src/eda - s550/plots/1. scatter with year.pdf")
 75  ggsave("RealEstate/src/eda - s550/plots/1. scatter with year.png")
 76
 77
 78  # facet scatter plots for 2020 by municipality
 79  re %>%
 80    dplyr::filter(year == 2020) %>%
 81    dplyr::select(-PIC, -year) %>%
 82    dplyr::filter(!is.na(total.assessment), !is.na(mill.rate)) %>%
 83    dplyr::group_by(municipality, tax.class) %>%
 84    dplyr::summarise(total.assessment = mean(total.assessment), mill.rate = mean(mill.rate)) %>%
 85    ggplot(aes(x = log(total.assessment), y = log(mill.rate))) +
 86    geom_point(color = viridis(20)[3]) +
 87    facet_wrap(.~tax.class, labeller = tax.classes) +
 88    labs(x = "log assessment value",
 89         y = "log mill rate") +
 90    theme(text = element_text(size = 18))
 91  ggsave("RealEstate/src/eda - s550/plots/2. scatter 2020 by municipality.pdf")
 92  ggsave("RealEstate/src/eda - s550/plots/2. scatter 2020 by municipality.png")
 93
 94
 95
 96  # facet line trends sample of 10
 97    re %>%
 98      dplyr::select(-PIC) %>%
 99      dplyr::filter(!is.na(total.assessment), !is.na(mill.rate)) %>%
100      dplyr::group_by(year, tax.class) %>%
101      dplyr::sample_n(size = 10) %>%
102      dplyr::ungroup() %>%
103      dplyr::group_by(year, tax.class, municipality) %>%
104      dplyr::summarise(total.assessment = mean(total.assessment), mill.rate = mean(mill.rate)) %>%
105      ggplot(aes(x = year, y = log(mill.rate), group = municipality)) +
106      geom_line(color = viridis(20)[3]) +
107      facet_wrap(.~tax.class, labeller = tax.classes) +
108      #scale_color_viridis_d() +
109      #theme(legend.position = "none") +
110      labs(x = "year",
111           y = "log mill rate") +
112      theme(text = element_text(size = 18),
113            axis.text.x = element_text(angle = 45, hjust = 1))
114    ggsave("RealEstate/src/eda - s550/plots/6. mill rate evolution sample.pdf")
115    ggsave("RealEstate/src/eda - s550/plots/6. mill rate evolution sample.png")
116
117      # facet line trends sample of 10 assessment
118    re %>%
119      dplyr::select(-PIC) %>%
120      dplyr::filter(!is.na(total.assessment), !is.na(mill.rate)) %>%
121      dplyr::group_by(year, tax.class) %>%
122      dplyr::sample_n(size = 10) %>%
123      dplyr::ungroup() %>%
124      dplyr::group_by(year, tax.class, municipality) %>%
125      dplyr::summarise(total.assessment = mean(total.assessment), mill.rate = mean(mill.rate)) %>%
126      ggplot(aes(x = year, y = log(total.assessment), group = municipality)) +
127      geom_line(color = viridis(20)[3]) +
128      facet_wrap(.~tax.class, labeller = tax.classes) +
129      #scale_color_viridis_d() +
130      #theme(legend.position = "none") +
131      labs(x = "year",
132           y = "log assessment value") +
133      theme(text = element_text(size = 18),
134            axis.text.x = element_text(angle = 45, hjust = 1))
135    ggsave("RealEstate/src/eda - s550/plots/6.1 assessment evolution sample.pdf")
136    ggsave("RealEstate/src/eda - s550/plots/6.1 assessment evolution sample.png")
137
138
139
140
```

```
141    # violin plots of mill rates accross tax classes, for 2020
142    re %>%
143      dplyr::filter(!is.na(total.assessment), !is.na(mill.rate), year == 2020) %>%
144      dplyr::group_by(PIC, tax.class) %>%
145      dplyr::summarise(total.assessment = mean(total.assessment), mill.rate = mean(mill.rate)) %>%
146      ggplot(aes(x = tax.class, y = log(mill.rate), fill = tax.class)) +
147      geom_violin(alpha = 0.5, width = 1) +
148      geom_boxplot(alpha = 0.75, width = 0.1) +
149      labs(x = "tax class",
150           y = "log mill rate") +
151      #scale_fill_viridis_d(begin=0, end=1) +
152      scale_fill_manual(values = c("#3E4A89FF", "#26828EFF","#B4DE2CFF" )) +
153      theme(legend.position = "none") +
154      scale_x_discrete(labels = c("01 - Residential", "05 - Industrial", "06 - Commercial")) +
155      theme(text = element_text(size = 18))
156    ggsave("RealEstate/src/eda - s550/plots/7. violin mill rates.pdf")
157    ggsave("RealEstate/src/eda - s550/plots/7. violin mill rates.png")
158
159
160    # violin plots of assessment values accross tax classes, for 2020
161    re %>%
162      dplyr::filter(!is.na(total.assessment), !is.na(mill.rate), year == 2020) %>%
163      dplyr::group_by(PIC, tax.class) %>%
164      dplyr::summarise(total.assessment = mean(total.assessment), mill.rate = mean(mill.rate)) %>%
165      ggplot(aes(x = tax.class, y = log(total.assessment), fill = tax.class)) +
166      geom_violin(alpha = 0.5, width = 1) +
167      geom_boxplot(alpha = 0.75, width = 0.1) +
168      labs(x = "tax class",
169           y = "log assessment value") +
170      #scale_fill_viridis_d(begin=0, end=1) +
171      scale_fill_manual(values = c("#3E4A89FF", "#26828EFF","#B4DE2CFF" )) +
172      theme(legend.position = "none") +
173      scale_x_discrete(labels = c("01 - Residential", "05 - Industrial", "06 - Commercial")) +
174      theme(text = element_text(size = 18))
175    ggsave("RealEstate/src/eda - s550/plots/7.1 violin assessment values.pdf")
176    ggsave("RealEstate/src/eda - s550/plots/7.1 violin assessment values.png")
```

Code 1: Code used for the exploratory data analysis.

# B   Modelling code

```
1    # preamble
2
3    suppressPackageStartupMessages(library(randomForest))
4    suppressPackageStartupMessages(library(tidyverse))
5    suppressPackageStartupMessages(library(lme4))
6    suppressPackageStartupMessages(library(ROCR))
7    suppressPackageStartupMessages(library(predictmeans))
8
9
10   ###########################
11   # data prep
12   dat <- readr::read_delim("test_train_data.txt", delim = ",", col_types = "ciccdddddc")
13   dat$municipality <- as.factor(dat$municipality)
14   dat$tax.class <- as.factor(dat$tax.class)
15
16   factors_tbl = dat %>%
17     group_by(municipality) %>%
18     count(name="mun_count", sort = TRUE) %>%
19     ungroup() %>%
20     mutate(perc = mun_count/sum(mun_count),
21            cum_perc = cumsum(perc)) %>%
22     arrange(desc(mun_count)) %>%
23     mutate(rank = row_number(),
24            municipality = fct_reorder(municipality, rank)) %>%
25     mutate(col_municipality = fct_collapse(municipality, other = levels(municipality)[-c(1:52)])) %>%
26     select(municipality, col_municipality)
27
28
```

```r
29 | med_assessment_by_municipality = dat %>%
30 |   left_join(factors_tbl) %>%
31 |   select(-c(municipality)) %>%
32 |   rename(municipality = col_municipality) %>%
33 |   filter(test.train == "train") %>%
34 |   group_by(municipality, tax.class, year) %>%
35 |   mutate(med_assessment = median(total.assessment, na.rm=TRUE)) %>%
36 |   select(municipality, year, tax.class, med_assessment) %>%
37 |   distinct(municipality, .keep_all = T)
38 |
39 | dat_mill <- dat %>%
40 |   left_join(factors_tbl) %>%
41 |   select(-c(municipality)) %>%
42 |   rename(municipality = col_municipality) %>%
43 |   left_join(med_assessment_by_municipality) %>%
44 |   group_by(PIC) %>%
45 |   mutate(next.assess = lead(med_assessment, order_by = year),
46 |          past.mill = lag(mill.rate, order_by = year)) %>%
47 |   arrange(PIC) %>%
48 |   group_by(municipality, year) %>%
49 |   mutate(n.prop = n()) %>%
50 |   arrange(desc(n.prop)) %>%
51 |   distinct(municipality, .keep_all = T) %>%
52 |   select(-c(tax, improvement.assessment, land.assessment, total.assessment))
53 |
54 | dat_as <- dat %>%
55 |   left_join(factors_tbl) %>%
56 |   select(-c(municipality)) %>%
57 |   rename(municipality = col_municipality) %>%
58 |   group_by(PIC) %>%
59 |   arrange(year) %>%
60 |   mutate(next.assess = lead(total.assessment, order_by = year),
61 |          past.mill = lag(mill.rate, order_by = year)) %>%
62 |   group_by(municipality) %>%
63 |   top_n(25, wt = total.assessment) %>%
64 |   arrange(municipality)
65 |
66 | train_mill <- dat_mill %>%  filter(test.train == "train")
67 | test_mill <- dat_mill %>%  filter(test.train == "test")
68 | train_as <- dat_as %>%  filter(test.train == "train")
69 | test_as <- dat_as %>%  filter(test.train == "test")
70 |
71 | ###################################
72 | # Random Forest
73 | set.seed(0)
74 |
75 | rf.mill <- randomForest(
76 |   mill.rate ~ tax.class + municipality + med_assessment + past.mill, na.action = na.omit, mtry = 4,
77 |   data=train_mill, ntree=500
78 | )
79 |
80 | save(rf.mill, file = "rf.mill.rda")
81 |
82 | #Evaluate variable importance
83 | importance(rf.mill)
84 | varImpPlot(rf.mill)
85 |
86 | rf.as <- randomForest(
87 |   next.assess ~ tax.class + municipality + total.assessment + mill.rate, na.action = na.omit,mtry =
88 |         4,
88 |   data=train_as
89 | )
90 |
91 | save(rf.as, file = "rf.as1.rda")
92 | importance(rf.as)
93 | varImpPlot(rf.as)
94 |
95 | yhat.bag <- predict(rf.mill,newdata=test_mill)
96 | plot(yhat.bag, test_mill$mill.rate, xlab="Predicted Mill Rate Using Test Set", ylab="Actual Mill Rate
      ")
97 | abline(0,1)
98 |
99 | yhat.bag1 <- predict(rf.as,newdata=test_as)
```

```
100  plot(yhat.bag1, test_as$next.assess, xlab="Predicted␣Asssessment␣Value␣Using␣Test␣Set", ylab="Actual␣
         Assessment␣Value", ylim=c(0, 10e8))
101  abline(0,1)
102
103  ################################
104  # Linear Mixed Effects Model
105
106  # both random slope and random intercept
107  # different rate of change of assessment value and mill rate as well as initial assessment value and
         mill rate for each municipality
108
109  lme.mill <- lmer(mill.rate ~ 1+ (1+year|municipality) + as.factor(year) + tax.class  + avg_assessment
         + past.mill, data = train_mill)
110
111  summary(lme.mill)
112
113  lme.as <- lmer(next.assess ~ 1+ (1+year|municipality) + as.factor(year) + tax.class + total.
         assessment + mill.rate, data = train_as)
114
115  summary(lme.as)
116
117  # check assumptions of lme
118
119  # Homogeneity of Variance
120  residplot(lme.mill)
121  residplot(lme.as)
```

Code 2: Code used for the modeling.

## C   Shiny app code

```
1
2   #global.R file
3
4   library(shiny)
5   library(readr)
6   library(ggplot2)
7   library(dplyr)
8   library(DT)
9   library(tidyr)
10  library(shinyjs)
11  library(randomForest)
12
13  #PIC used for some testing
14  # CA-BC-200-001019632060000
15
16
17  # for bookmarking button
18  enableBookmarking("url")
19
20  # read data
21  dat <- readr::read_delim("test_train_data.txt",
22                           delim = ",", col_types = "ciccdddddf")
23
24  # creates dataset with only the top 52 municipalities
25  counts <- dat %>%
26    count(municipality, sort = TRUE)
27
28  datshort <- dat %>%
29    filter(municipality %in% counts$municipality[1:52])
30
31
32  # load dataset used for rf.mill
33  rfdat <- readRDS("rf_data.rds")
34
35  # create dataset used for rf.as
36  asdat <- readRDS("as_data.rds")
37
38  ################################
```

```
39   #################################
40   # ui.R file
41
42
43   # Mallory - STAT 550 2020###
44   # This is the user interface version of the shiny app
45
46   library(shinythemes)
47   library(png)
48
49   ui <- fluidPage(theme = shinytheme("cerulean"), #maybe journal theme?
50
51                    # header
52                    div(id = "headerSection",
53                        h2("BC Mill Rate & Assessment Value Predictions"),
54
55                        span(
56                          style = "font-size: 1em",
57                          # authors
58                          span("Created by "),
59                          a("Gian Carlo Diluvi, Vittorio Romaniello, Sophia Li & Mallory Flynn",
60                            href = "https://www.stat.ubc.ca"),
61                          HTML("&bull;"),
62                          # date
63                          span("April 2020"),
64                          HTML("&bull;"),
65                          # Shiny app code link
66                          span("Code"),
67                          a("on GitHub",
68                            href = "https://github.com/STAT450-550/RealEstate/tree/master/src/shiny_app")
69                          )
70                    ),
71                    br(),
72                    br(),
73
74                    # all content goes here, and is hidden initially until the page fully loads
75                    sidebarLayout(
76                      sidebarPanel(
77                        # tabsetPanel(
78                        #   tabPanel("User Inputs",
79
80                        # Only show the following for assessment predictions:
81                        # Use PIC?
82                        checkboxInput("picInput", "Use PIC?", value = FALSE),
83                        selectInput("typeInput", "Estimate Type",
84                                    c("Select", "Assessment Value", "Mill Rate"),
85                                    selected = "Select"),
86
87                        # If using PIC:
88                        conditionalPanel("input.picInput",
89                                         textInput("identInput", "PIC:", placeholder = NULL)),
90
91                        # If PIC is not available:
92                        conditionalPanel("!input.picInput",
93
94                                         # for municipality
95                                         selectInput("municipalityInput", "Municipality:",
96                                                     c("-",sort(unique(datshort$municipality))),
97                                                     selected = "-"),
98
99                                         # for Tax Class code
100                                        selectInput("taxclassInput", "Tax Class Code:",
101                                                    c("-",sort(unique(dat$tax.class))),
102                                                    selected = "-"),
103
104                                        #conditional input for estimate type
105                                        conditionalPanel("input.typeInput == 'Assessment Value'",
106                                                         numericInput("assessmentInput",
107                                                                      "Current Assessment Value:",
108                                                                      value = 70000000,
109                                                                      min = 4241700,
110                                                                      max = 10000000000))
111                                        ),
```

13

```
112
113
114
115                             # button to update the data
116                             shiny::hr(),
117                             actionButton("updateButton", "Update"),
118
119
120
121                             # source of data as a footer - Altus Group image not loading
122                             br(),
123                             br(),
124                             p("Generated␣using␣data␣from␣",
125                               a("the␣Altus␣Group␣Ltd.",
126                                   href = "https://www.altusgroup.com",
127                                   target = "_blank")),
128                             a(img(src = "altusgroupimg.png", alt = "Altus␣Group",
129                                     height = 63, width = 150),
130                               href = "https://www.altusgroup.com",
131                               target = "_blank"),
132                             br(),
133                             br(),
134                             br(),
135                             br(),
136                             bookmarkButton()
137                             ),
138
139
140                         # main panel with Estimate tab and plot tab with mill rates
141                         # or assessment values over time
142                         mainPanel(h4(textOutput("resultsText")),
143                                 tabsetPanel(
144                                   tabPanel("Estimate",
145                                           br(),
146                                           verbatimTextOutput("results")),
147                                   tabPanel("Plot",
148                                             br(),
149                                             plotOutput("coolplot"))
150                                 )
151                     )
152                 )
153                 )
154
155
156 #########################################
157 #########################################
158 # server.R file
159
160 # Mallory - STAT 550 2020###
161 # This is the server file of the shiny app
162
163 # fix main title when PIC is checked but empty
164 # load rdas for each so that estimates can be made
165 # fix select input to choose only the top 52 categories and other
166
167
168 # in case modified data needs to be accessed
169 source("helpers.R")
170
171 # load models - RF for mill rate predictions and for assessment value predictions
172 load("rf.mill.rda")
173 load("rf.as.rda")
174
175
176 # server:
177 server <- function(input, output, session) {
178
179   filtered <- reactive({
180
181     # Update when following inputs are changed
182     input$updateButton
183
184     newdata <- datshort
185     d <- NULL
```

```
186        #print(dim(dat))
187
188        # Filter data based on the user inputs
189        isolate({
190              # If using PIC:
191              if(input$picInput && input$identInput!=""){
192                d <- newdata  %>%
193                  filter(PIC == input$identInput)
194              }
195
196
197              # If not using PIC, filter by municipality and tax class:
198              if(!input$picInput){
199
200                d <- newdata %>%
201                  filter(tax.class == input$taxclassInput,
202                         municipality == input$municipalityInput)
203              }
204
205          })
206
207        # return filtered data
208        if(dim(d)[1]==0){
209          d <- NULL
210        }
211        d
212
213        })
214
215        ########## PLOTTING TAB ###################
216        # Add plots of either mill rate or assessment value over time to plot tab
217
218      # create mill rate plot that reacts to inputs
219      millrateplot <- reactive({
220        input$updateButton
221
222        data <- filtered()
223
224
225        isolate({
226        if(is.null(data)){
227          p <- paste("No corresponding data to plot.")
228        }
229
230        # plot mill rates over time for municipality chosen for mill rate predicitons
231        if (input$typeInput == 'Mill Rate'){
232          p <- ggplot(data, aes(x = year, y = mill.rate)) +
233            geom_line(color="#FF3333") +
234            geom_point(color="#FF3333") +
235            theme_minimal() +
236            xlab("Year") +
237            ylab("Mill Rate") +
238            ggtitle("Municipal Mill Rate Over Time")
239        }
240        })
241
242        p
243
244    })
245
246      # create assessment plots that react to user inputs
247      assessplot <- reactive({
248        input$updateButton
249
250        data <- filtered()
251
252        isolate({
253          if(is.null(data)){
254            p <- paste("No corresponding data to plot.")
255          }
256
257        # plot assessment values over time
258          if(input$picInput && input$identInput!=""){
259            if(input$typeInput == 'Assessment Value') {
```

15

```
260                 #print("ggplotting assessment values")
261                 p <- ggplot(data, aes(x = year, y = total.assessment)) +
262                   geom_line(color = "#56B4E9") +
263                   geom_point(color = "#56B4E9") +
264                   theme_minimal() +
265                   xlab("Year") +
266                   ylab("Assessment Value") +
267                   ggtitle("Assessment Values Over Time")
268             }
269         }
270         else{
271             p <- paste("No corresponding data to plot.")
272         }
273         })
274
275      p
276
277    })
278
279    # output one of the above plots onto UI
280     output$coolplot <- renderPlot({
281       if (input$typeInput != 'Select'){
282         if(input$typeInput == 'Assessment Value'){
283             assessplot()
284         }
285
286         else{
287             millrateplot()
288         }
289         }
290
291       else{
292         return()
293       }
294         })
295
296
297    ###### ESTIMATE TAB ####################
298    # give predictions given user inputs for mill rate or assessment value
299    estimates <- reactive({
300      input$updateButton
301
302      isolate({
303         if(is.null(filtered())){
304             pred <- paste("No data.")
305         }
306
307         else{
308           # If using PIC:
309
310           # If doing Mill Rate prediciton:
311           if(input$typeInput == 'Mill Rate'){
312             #print("doing mill rate prediction")
313
314             # extract latest mill rate
315             past20 <- filtered() %>%
316               filter(year == 2020)
317
318             # if 2020 column is empty, it will break by condition on mean mill rate=0
319             meanmillrate <- mean(na.omit(past20$mill.rate))
320             print(meanmillrate)
321
322             # put data together in the way rfmill expects as input call it inputdata
323             # columns include tax.class, municipality, total.assessment, past.mill
324             meanassess <- mean(na.omit(past20$total.assessment))
325             print(meanassess)
326
327             pred.data <- cbind(filtered()$tax.class[1],
328                               filtered()$municipality[1],
329                               meanassess,
330                               meanmillrate)
331             pred.data <- as.data.frame(pred.data, stringsAsFactors = FALSE)
332             colnames(pred.data) <- c('tax.class', 'municipality',
333                                     'avg_assessment','past.mill')
```

16

```
334
335
336             pred.data$past.mill <- as.numeric(pred.data$past.mill)
337             pred.data$avg_assessment <- as.numeric(pred.data$avg_assessment)
338             pred.data$municipality <- factor(pred.data$municipality,
339                                          levels = levels(rfdat$municipality))
340             pred.data$tax.class <- factor(pred.data$tax.class,
341                                      levels = levels(rfdat$tax.class))
342             print(pred.data)
343
344             # predict next mill rate using random forest
345             pred <- round(predict(rf.mill, newdata = pred.data), 2)
346
347             if(meanmillrate==0 || is.na(pred)){
348               pred <- paste("No previous mill rate found in data.")
349             }
350
351             pred
352
353             }
354
355
356          # If doing Assessment Value prediction:
357          if(input$typeInput == 'Assessment Value'){
358
359             print("doing assessment value prediction")
360
361             # extract latest assessment value
362             past20 <- filtered() %>%
363               filter(year == 2020)
364
365             # if using PIC and 2020 column is NA for this property's
366             # assessment value, it will break
367             if(input$picInput){
368               print("using PIC")
369               if(input$identInput != ""){
370                 if(!is.na(past20$total.assessment)){
371                   last.assess <- past20$total.assessment
372                   }
373                 else{
374                   last.assess <- 0
375                 }
376               }
377               print(last.assess)
378             }
379
380             # if not using PIC, assessment value must come from user input
381             else{
382               print("not using PIC")
383               if(input$assessmentInput != ""){
384                 last.assess <- input$assessmentInput
385                 }
386               else{
387                 last.assess <- 0
388                 }
389               print(last.assess)
390             }
391
392
393             # put data together in the way rfmill expects as input call it inputdata
394             # columns include tax.class, municipality, total.assessment, and mill.rate
395             print(head(filtered()))
396             pred.data <- cbind(filtered()$tax.class[1],
397                                filtered()$municipality[1],
398                                last.assess,
399                                past20$mill.rate[1])
400
401             pred.data <- as.data.frame(pred.data, stringsAsFactors = FALSE)
402             colnames(pred.data) <- c('tax.class', 'municipality',
403                                      'total.assessment', 'mill.rate')
404
405
406             pred.data$mill.rate <- as.numeric(pred.data$mill.rate)
407             pred.data$total.assessment <- as.numeric(pred.data$total.assessment)
```

```
408              pred.data$municipality <- factor(pred.data$municipality,
409                                          levels = levels(asdat$municipality))
410              pred.data$tax.class <- factor(pred.data$tax.class,
411                                          levels = levels(asdat$tax.class))
412              print(pred.data)
413
414              # predict next assessment value using random forest
415              pred <- round(predict(rf.as, newdata = pred.data),2)
416
417              if(last.assess==0 || is.na(pred)){
418                pred <- paste("Missing required data.")
419              }
420            }
421          }
422          })
423
424
425      pred
426      print(pred)
427
428      })
429
430    # create estimates as text for output
431    estimatestext <- reactive({
432      input$updateButton
433
434
435      # If using PIC:
436      if(input$picInput){
437        if(input$identInput!=""){
438          if(input$typeInput == 'Mill Rate'){   #need to be changed to extract values
439            return(paste("Mill rate prediction for class",
440                         filtered()$tax.class[1],
441                         "in", filtered()$municipality[1], "- \n",
442                         estimates(), sep = " "))  #  RETURN PREDICTION
443          }
444
445          if(input$typeInput == 'Assessment Value'){
446            return(paste("Predicted next assessment value of property \n", input$identInput,
447                         "-", estimates(), sep = " "))  # RETURN PREDICTION
448          }
449
450          if(input$typeInput == 'Select'){
451            return("Enter prediction type.")
452          }
453          }
454
455        else{
456          return("Please enter PIC.")
457        }
458        }
459
460      # If not using PIC:
461      else{
462        if(input$typeInput == 'Mill Rate'){
463          return(paste("Mill rate prediction for class", input$taxclassInput,
464                       "in", input$municipalityInput, "- \n",
465                       estimates(), sep = " ")) # RETURN PREDICTION
466        }
467
468        if(input$typeInput == 'Assessment Value'){
469
470          if(!is.null(filtered())){
471            md <- asdat %>%
472              filter(municipality == input$municipalityInput)
473
474            minm <- min(na.omit(md$total.assessment))
475            print(minm)
476
477            maxm <- max(na.omit(md$total.assessment))
478            print(maxm)
479          }
480
481          if(input$municipalityInput == '-' ||
```

```
482              input$taxclassInput == '-'){
483           return(paste("Complete␣user␣inputs."))
484         }
485
486          return(paste("Predicted␣next␣assessment␣value␣-␣\n",
487                    estimates(), "\n␣Valid␣prediction␣range␣for␣this␣municipality␣is",
488                    minm, "-", maxm, sep = "␣"))  # RETURN PREDICTION
489       }
490
491       if(input$typeInput == 'Select'){
492          return("Enter␣prediction␣type.")
493       }
494       }
495    })
496
497   # output the estimates text in the main panel
498   output$results <- renderText({
499     estimatestext()
500     })
501
502
503   # Titles text for main panel title - describes prediction type or PIC
504   # number if applicable
505   titles <- reactive({
506     input$updateButton
507
508     data <- filtered()
509
510     if(!input$picInput && input$typeInput == 'Select' && input$municipalityInput == '-' &&
511        input$taxclassInput == '-'){
512       return(paste(""))
513     }
514
515     if(is.null(data)){
516        return(paste("Could␣not␣find␣matching␣data."))
517     }
518
519     else{
520       if(input$picInput){
521         if(input$identInput !=""){
522
523            if(input$typeInput == 'Assessment␣Value'){
524              return(paste("Assessment␣value␣for", input$identInput, sep = "␣")) #ADD PREDICTION HERE
525            }
526
527            if(input$typeInput == 'Mill␣Rate'){
528              return(paste("Class", filtered()$tax.class[1],
529                        "mill␣rate␣for", filtered()$municipality[1], sep = "␣")) #ADD PREDICTION
                              HERE
530            }
531
532            if(input$typeInput == 'Select'){
533              return(paste("Select␣prediction␣type."))
534            }
535         }
536
537         else{
538           return(paste(""))
539         }
540       }
541
542       else{
543           if(input$typeInput == 'Mill␣Rate'){
544             return(paste("Class", input$taxclassInput,
545                        "mill␣rate␣for",
546                        input$municipalityInput, sep = "␣"))
547         }
548
549         else{
550           if(input$typeInput == 'Select'){
551             return(paste("Select␣prediction␣type."))
552           }
553           else{
554             return(paste("Assessment␣Value␣prediction␣for␣class",
```

19

```
555                                    input$taxclassInput, "property \n in",
556                                    input$municipalityInput, sep = " "))
557                  }
558               }
559            }
560         }
561     })
562
563
564     output$resultsText <- renderText({
565       titles()
566     })
567  }
568
569  #######################################
570  #######################################
571  # helpers.R file
572  # This file will have helpers for the model and the loading of data for models
573
574  # function to modify data for random forest (if needed)
575
576  # used to create mill rate data for random forest
577  # no longer needed for shiny app; included for future use if needed
578  rfData <- function(data) {
579
580    dat <- data
581    dat$municipality <- as.factor(dat$municipality)
582    dat$tax.class <- as.factor(dat$tax.class)
583
584    factors_tbl = dat %>%
585      group_by(municipality) %>%
586      count(name="mun_count", sort = TRUE) %>%
587      ungroup() %>%
588      mutate(perc = mun_count/sum(mun_count),
589             cum_perc = cumsum(perc)) %>%
590      arrange(desc(mun_count)) %>%
591      mutate(rank = row_number(),
592             municipality = fct_reorder(municipality, rank)) %>%
593      mutate(col_municipality = fct_collapse(municipality, other = levels(municipality)[-c(1:52)])) %>%
594      select(municipality, col_municipality)
595
596
597    avg_assessment_by_municipality = dat %>%
598      left_join(factors_tbl) %>%
599      select(-c(municipality)) %>%
600      rename(municipality = col_municipality) %>%
601      filter(test.train == "train") %>%
602      group_by(municipality, tax.class, year) %>%
603      mutate(avg_assessment = mean(total.assessment, na.rm=TRUE)) %>%
604      select(municipality, year, tax.class, avg_assessment) %>%
605      distinct(municipality, .keep_all = T)
606
607    dat_mill <- dat %>%
608      left_join(factors_tbl) %>%
609      select(-c(municipality)) %>%
610      rename(municipality = col_municipality) %>%
611      left_join(avg_assessment_by_municipality) %>%
612      group_by(PIC) %>%
613      mutate(next.assess = lead(avg_assessment, order_by = year),
614             past.mill = lag(mill.rate, order_by = year)) %>%
615      arrange(PIC) %>%
616      group_by(municipality, year) %>%
617      mutate(n.prop = n()) %>%
618      arrange(desc(n.prop)) %>%
619      distinct(municipality, .keep_all = T) %>%
620      select(-c(tax, improvement.assessment, land.assessment, total.assessment))
621
622    dat_mill
623  }
624
625  # used to create assessment value data for random forest
626  # no longer needed for shiny app; included for future use if needed
627  asData <- function(data) {
628
```

```
629    dat <- data
630    dat$municipality <- as.factor(dat$municipality)
631    dat$tax.class <- as.factor(dat$tax.class)
632
633    factors_tbl = dat %>%
634      group_by(municipality) %>%
635      count(name="mun_count", sort = TRUE) %>%
636      ungroup() %>%
637      mutate(perc = mun_count/sum(mun_count),
638             cum_perc = cumsum(perc)) %>%
639      arrange(desc(mun_count)) %>%
640      mutate(rank = row_number(),
641             municipality = fct_reorder(municipality, rank)) %>%
642      mutate(col_municipality = fct_collapse(municipality, other = levels(municipality)[-c(1:52)])) %>%
643      select(municipality, col_municipality)
644
645
646    avg_assessment_by_municipality = dat %>%
647      left_join(factors_tbl) %>%
648      select(-c(municipality)) %>%
649      rename(municipality = col_municipality) %>%
650      filter(test.train == "train") %>%
651      group_by(municipality, tax.class, year) %>%
652      mutate(avg_assessment = mean(total.assessment, na.rm=TRUE)) %>%
653      select(municipality, year, tax.class, avg_assessment) %>%
654      distinct(municipality, .keep_all = T)
655
656    dat_as <- dat %>%
657      left_join(factors_tbl) %>%
658      select(-c(municipality)) %>%
659      rename(municipality = col_municipality) %>%
660      group_by(PIC) %>%
661      arrange(year) %>%
662      mutate(next.assess = lead(total.assessment, order_by = year),
663             past.mill = lag(mill.rate, order_by = year)) %>%
664      group_by(municipality) %>%
665      top_n(25, wt = total.assessment) %>%
666      arrange(municipality)
667
668    dat_as
669 }
```

Code 3: Code used for the Shiny app.