

Biostatistics 625 Final Project Report

Group 2 members: Mukai Wang, Kangping Yang, Litian Zhou

December 22, 2019

Introduction

Housing costs in China have witnessed turbulent fluctuation in the past few years, especially in the country's capital Beijing. In this project, we design and implement an R shiny app that offers interactive data visualization and statistical insights about the housing price trend in Beijing. The data source is from Kaggle. It contains 300k+ transaction logs between 2011 and 2017.

Dataset Summary and Storage

The raw data contains 26 variables. After we remove columns that are not useful (e.g. ids) and columns with too many missing values, we retain 20 variables that we are interested in. For details of the variables, please refer to our codebook in the github repo.

In order to better explain the association between price and time, we decide to generate a new variable called **season** based on the tradetime. The **season** variable sets the first season of 2010 as 1, the second season of 2010 as 2, etc until the first season of 2018 as 33.

Because we need to subset the data in a variety of ways efficiently for the shiny app to display the data from various aspects, we decide to store the data in a relational database. Because we plan to deploy our shiny app online, an online database is our desired option. We end up using Amazon Relational Database Service. Because our data is not too large, a free tier version is sufficient. The following example illustrates how we can query data from the database. For more details, please refer to our github repo.

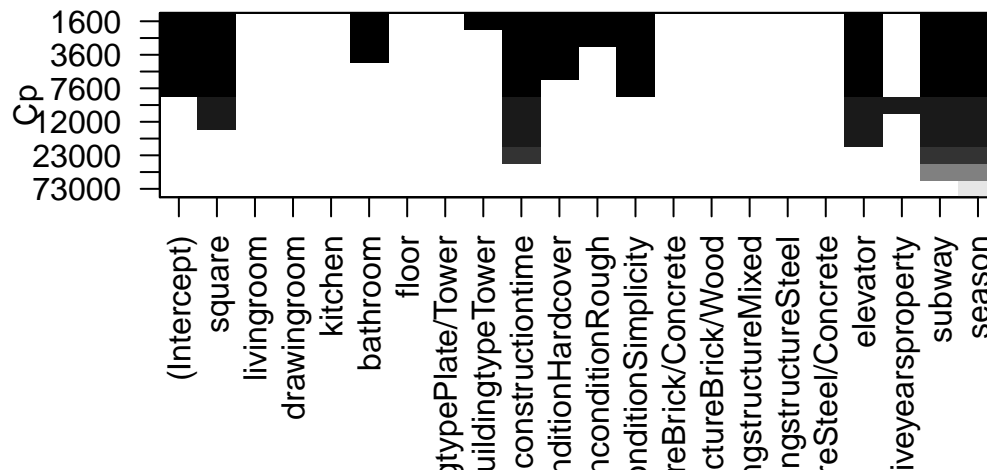
```
# data_query_example.R gives us the handler to query data from the database
source("../data_cleaning_storage/data_query_example.R")
# gen_data is a function that acts as an easy api for
# querying data based on selected district and building type
# gen_full_data is a similar function that
# returns data with all covariates above to demonstrate full model.
exampledata = gen_full_data("all", c())
```

Statistical Model

We try to see if a linear model can capture the trend between price and other covariates. To decide which covariates need to be taken into account, we start from the full model with all the covariates listed in the codebook.

```
library(leaps)
full_model <- lm(price ~ square + livingroom + drawingroom + kitchen + bathroom + floor +
  buildingtype + constructiontime + renovationcondition +
  buildingstructure + elevator + fiveyearsproperty +
  subway + season, data = exampledata$data)
mat_full <- model.matrix(full_model)
selection <- regsubsets(mat_full, exampledata$data$price, int = F, really.big = T,
  method = "exhaust", nvmax = 11)
```

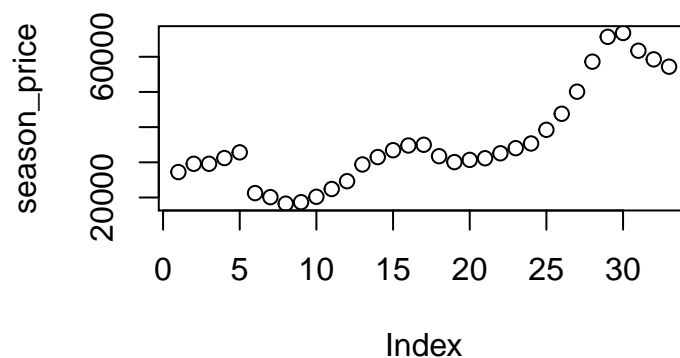
```
plot(selection, scale = "Cp")
```



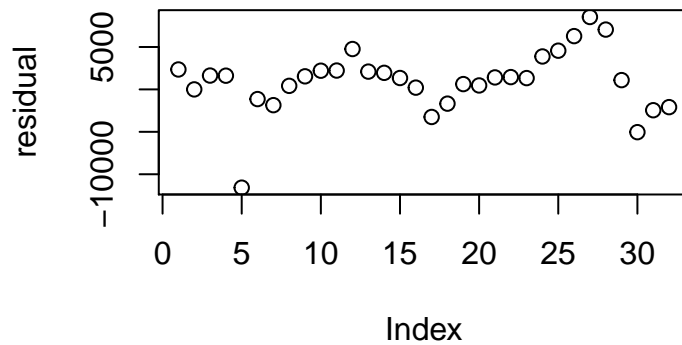
Based on the Mallows's Cp statistic, We choose covariates square,bathroom,buildingtype,renovationcondition, district, elevator, subway and season. To increase the model accuracy we decide to treat season as a factor.

```
exampledata$data$season <- as.factor(exampledata$data$season)
model_select <- lm(price ~ square + bathroom + buildingtype + constructiontime +
                    renovationcondition + elevator + subway + district +
                    season, data = exampledata$data)
season_price <- c(0,model_select$coefficients[23:54])
district_stat <- table(exampledata$data$district)
weight_dist <- district_stat/sum(district_stat)
reference <- model_select$coefficients[1] +
             mean(exampledata$data$bathroom)*model_select$coefficients[3] +
             sum(weight_dist*c(0,model_select$coefficients[12:22]))
```

```
season_price <- season_price + reference
##extract season coefficients as average selling price for time-series references.
plot(season_price) ## see the trend of housing price across seasons
```

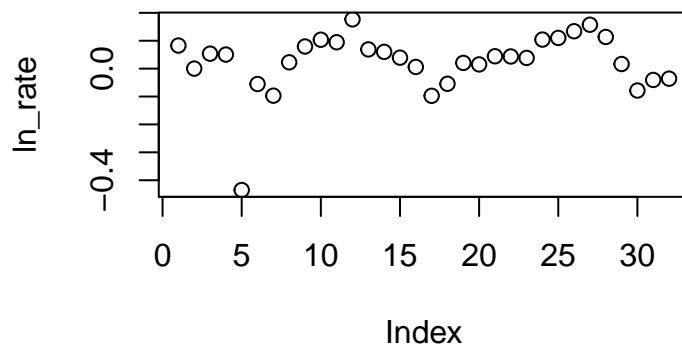


```
residual <- season_price[2:length(season_price)]-season_price[-length(season_price)]
plot(residual) ## residual plot between consecutive seasons
```

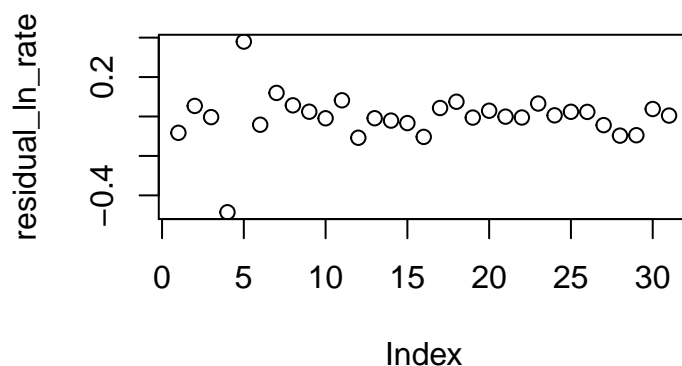


We find strong correlation in terms of price differences between adjacent seasons. They tend to become unstable as the price goes up. So we decide to calculate the log of the rate of the price change and observe its residual.

```
price_rate <- season_price[2:length(season_price)]/season_price[-length(season_price)]
ln_rate <- log(price_rate)
plot(ln_rate)
```



```
residual_ln_rate <- ln_rate[2:length(ln_rate)]-ln_rate[-length(ln_rate)]
plot(residual_ln_rate)
```



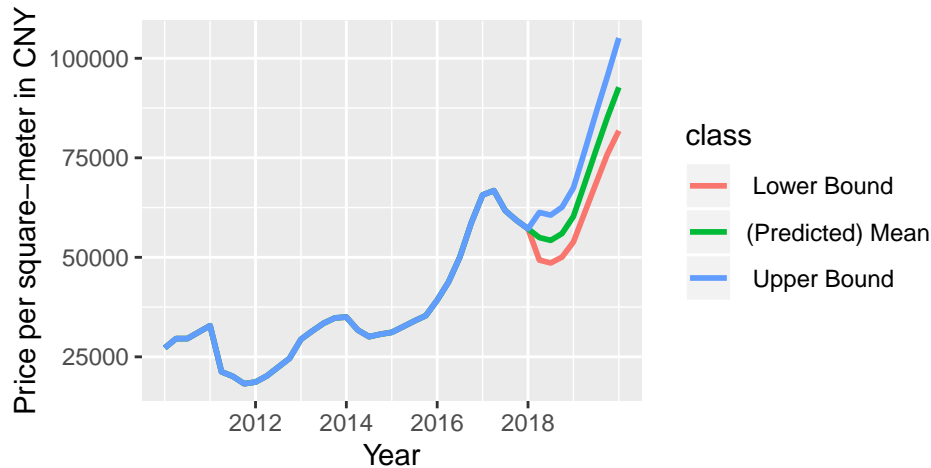
The residual plot looks better as the log rate change is closer to 0 and tends to become stable across the time. Therefore, we assume the log rate difference between adjacent seasons tend to be a constant around zero. For time series model $ARIMA(p, d, q)$, this assumption indicates $d = 1$. As for p , we found the price tend to follow a similar pattern in about every two-year period. So we decide to choose $p = 8$ (8 seasons i.e. 2 years). As for q , we found the residual of log rate tend to be auto-regressive by 1 season. We decide to choose $q = 1$. The prediction model we use is thus $ARIMA(8,1,1)$.

```
# the model fitting function is in another script
source("../Prediction_Model_Function.R")
# the function accepts the object returned from the gen_data function.
```

```
modeloutput = Filter_model(exampleddata)
```

The predicted house price in 2018 to 2019 from the fitted ARIMA(8,1,1) is shown in the following price trend graph:

```
plots = ggplot(modeloutput$beta_data, aes(x = year, y = price , color = class)) +
  geom_line(size=1) + labs(x = 'Year' ,y = 'Price per square-meter in CNY') +
  scale_x_continuous(breaks = c(2012,2014,2016,2018))
plots
```



App Interface

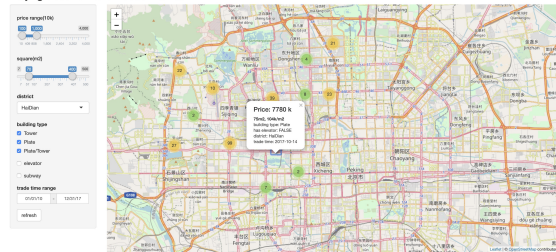
We integrate data visualization and statistical inference into one webpage supported by Rshiny. The interface has four panels: map & markers panel, filter panel, plots panel and statistical inference panel.

In the left panel, user can apply filters to the markers shown on the map, including Price Range, house area (square), district, building type, if houses are with a elevator or close to subway stations, also the trading time range. The main panel shows a map of Beijing, with no markers in the beginning. By applying the filters and hit “refresh,” the map will random sampling seven hundred records which satisfies the filters for the sake of reactivity.

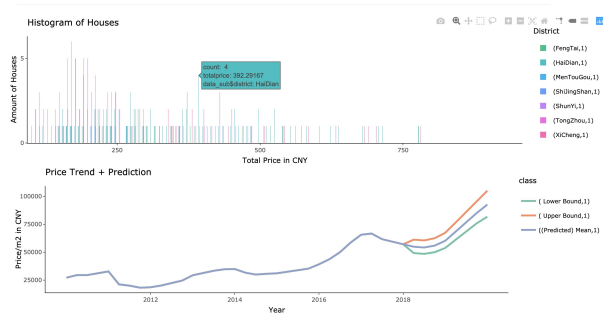
To make the map look cleaner, we cluster the nearby points into a group with the amount number on it. Users can click the cluster to expand the markers. User can also click each marker and get pop-up trade informations including price, housing area and trade time.

In the lower panel, there are two plots and statistical inference drawn from the user-specified subset. Please note that histogram uses the same filter as the left panel, while the “Price Trend + Prediction” plot uses all the data in the specified district and building type. The reason is that a wider range of records can bring better prediction and more meaningful interpretation.

Beijing Second-hand House Market



Plots:



Deliverable and Code Repository

Shiny App: https://biostat-umich.shinyapps.io/Beijing_Housing_Price/

GitHub: https://github.com/LitianZhou/Beijing_Housing_Price

Individual Contribution

Main tasks:

Shinyapp interface and visualization: Litian Zhou;

Database construction, maintainance and data query: Mukai Wang;

Model building and variable selection: Kangping Yang.

All the works are done as a group, so each group member contributes to all other tasks as well.

Conclusion

In this project, We manage to set up an R shiny app that produces interactive visualizations and statistical inference for Beijing housing prices between 2010 and 2018. The app can serve as a useful tool to review housing price data in the past and it can also give some predictions into the near future. All the team members enhance their technical abilities (including R Shiny, Amazon Database and time series analysis) through implementing this app. Overall the project is successful and meets the original expectation.