

# Predicting House Prices

## Quantitative Methods 1 (Tutorial)

Conor, Linette, Lucas, Minh

### 1 The intuition behind variable selection

To predict the Adjusted Sale Price of houses, we considered three major factors:

- the neighbourhood;
- the size of the house; and
- the quality of the construction.

Figure 1 below summarizes the variables we selected to represent each of those characteristics and their relationship.

```
1 pairs(train[c(9,12,16,22,24)])
```

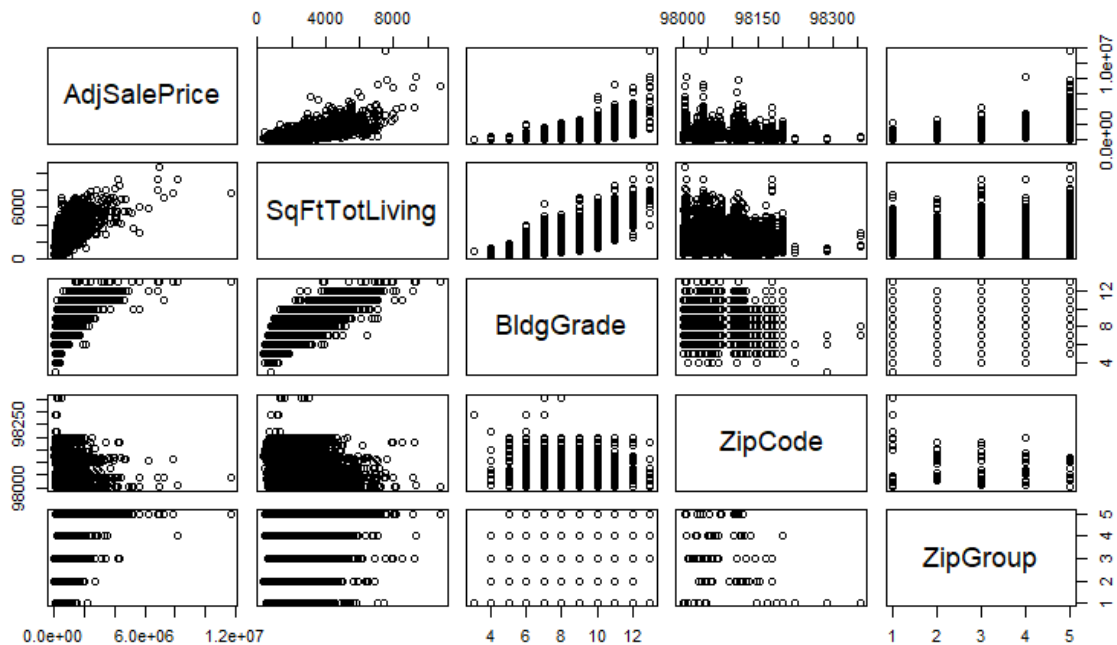


Figure 1: Associations of `AdjSalePrice` with `SqFtTotLiving`, `BldgGrade`, and `ZipCode`

ZipCode represents **neighborhood**. We grouped zip codes based on house prices in the variable ZipGroup to reflect the socioeconomic profile of each neighborhood.

```

1 zip_group_pr <- train %>%
2   group_by(ZipCode) %>%
3   summarise(med_price = median(AdjSalePrice),
4             count = n()) %>%
5   arrange(med_price) %>%
6   mutate(cumul_count = cumsum(count),
7          ZipGroup = ntile(cumul_count, 5))
8
9 train <- train %>%
10  left_join(select(zip_group_pr, ZipCode, ZipGroup), by = "ZipCode")

```

SqFtTotLiving represents the **size of the house**. We avoided adding to the model the variables SqFtLot, SqFtFinBasement, NbrLivingUnits Bathrooms, and Bedrooms because they are also related to the size of the construction. That decision allowed us to reach a more **parsimonious** model, preventing us from losing degrees of freedom and from incurring in issues related to **collinearity**. The risk of collinearity among those variables is visible in Figure 2 below:

```

1 pairs(train[c(12,11,13,14,15)])

```

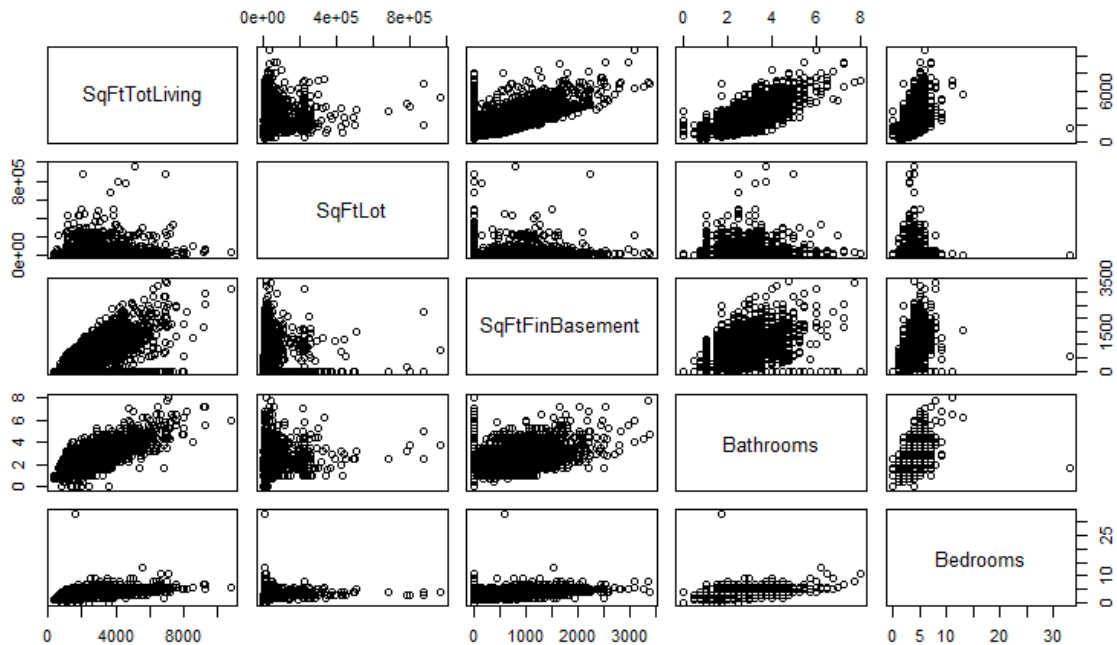


Figure 2: Risk of collinearity among SqFtLot, SqFtFinBasement, Bathrooms, and Bedrooms

Lastly, BldgGrade represents the **quality of the construction**. For the same reasons stated above (parsimony, preventing the loss of degrees of freedom, and avoiding collinear-

ity issues), we decided not to add `YrBuilt`, `YrRenovated`, and `NewConstruction`, which are also related to the quality of the construction.

## 2 The intuition behind the interaction effect

To complete the model, we added an interaction effect between `ZipGroup` and `SqFtTotLiving`. The interaction reflects the idea that each additional square foot in the building will affect house prices differently depending on the neighborhood. As a consequence, the slope of `SqFtTotLiving` on `AdjSalePrice` in expensive neighborhoods will be steeper than in popular ones (as can be seen in Figure 3).

```
1 ggplot(train, aes(SqFtTotLiving, AdjSalePrice, group = ZipGroup)) +  
2   geom_point(aes(colour = ZipGroup)) +  
3   geom_smooth(method = "lm", aes(colour = ZipGroup))
```

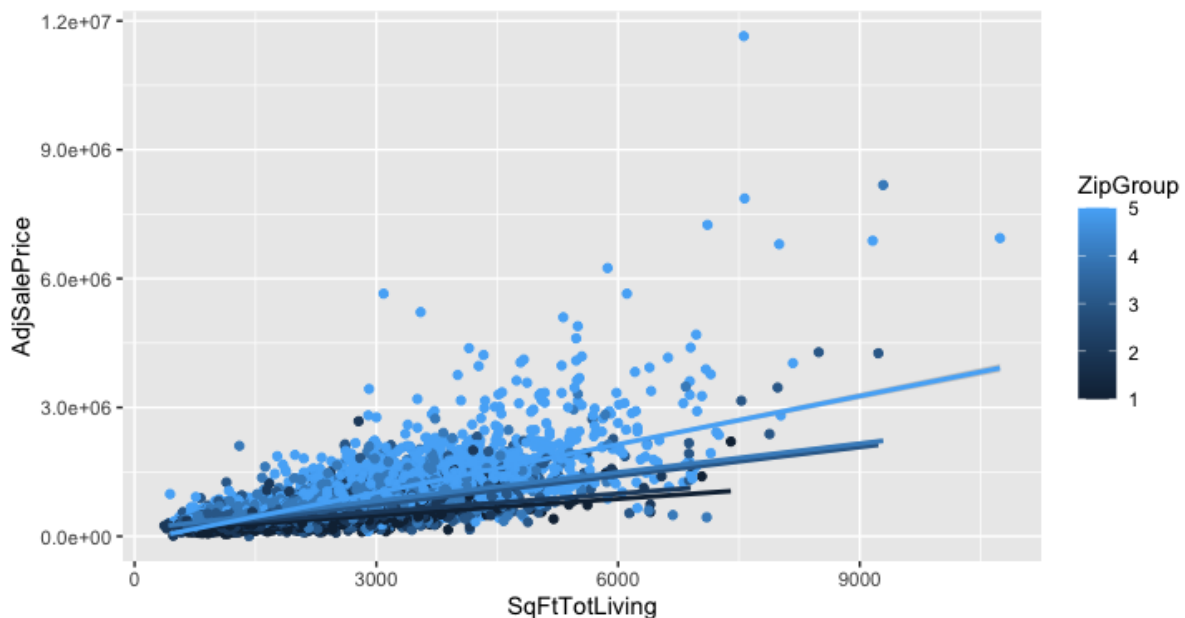


Figure 3: The effect of `SqFtTotLiving` on `AdjSalePrice` by `ZipGroup`

## 3 Final model

The effects of `SqFtTotLiving`, `BldgGrade`, `as.factor(ZipGroup)`, and `SqFtTotLiving:ZipGroup` on `AdjSalePrice`

```
1 mod1 <- lm(AdjSalePrice ~ SqFtTotLiving, data = train)  
2 mod2 <- lm(AdjSalePrice ~ BldgGrade, data = train)  
3 mod3 <- lm(AdjSalePrice ~ ZipGroup, data = train)  
4 mod4 <- lm(AdjSalePrice ~ SqFtTotLiving + BldgGrade + as.factor(ZipGroup) +  
5   SqFtTotLiving:ZipGroup,  
6   data = train,  
   na.action = na.omit)
```

7 `stargazer(mod1, mod2, mod3, mod4)`

Table 1: The effects of SqFtTotLiving, BldgGrade, `as.factor(ZipGroup)`, and `SqFtTotLiving:ZipGroup` on `AdjSalePrice`

	<i>Dependent variable:</i>			
	AdjSalePrice			
	(1)	(2)	(3)	(4)
SqFtTotLiving	294.357*** (2.132)			−45.169*** (5.890)
BldgGrade		221,513.500*** (1,695.289)		74,339.440*** (2,324.481)
ZipGroup			134,825.800*** (1,767.656)	
<code>as.factor(ZipGroup)2</code>				−60,600.140*** (6,261.020)
<code>as.factor(ZipGroup)3</code>				−112,855.500*** (7,949.155)
<code>as.factor(ZipGroup)4</code>				−168,392.000*** (9,795.236)
<code>as.factor(ZipGroup)5</code>				−257,812.900*** (13,626.880)
<code>SqFtTotLiving:ZipGroup</code>				64.755*** (1.442)
Constant	−47,126.100*** (4,843.068)	−1,136,579.000*** (13,173.570)	138,024.100*** (6,085.377)	−232,072.400*** (15,604.750)
Observations	20,340	20,340	20,340	20,340
R <sup>2</sup>	0.484	0.456	0.222	0.622
Adjusted R <sup>2</sup>	0.484	0.456	0.222	0.621
Residual Std. Error	278,257.400 (df = 20338)	285,528.100 (df = 20338)	341,480.700 (df = 20338)	238,266.300 (df = 20332)
F Statistic	19,053.750*** (df = 1; 20338)	17,073.130*** (df = 1; 20338)	5,817.690*** (df = 1; 20338)	4,770.375*** (df = 7; 20332)

Note:

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01