

# A Multiple Linear Model for Toronto and Mississauga House Prices

Tanyilan, Id 1004701548

December 5, 2020

## I. Data Wrangling

### 1.1 Sampling Data

We randomly selected 150 cases from the raw data. The id of the cases are

##	[1]	51	99	114	131	60	158	17	133	30	92	145	142	150	148	118	117	4	115
##	[19]	64	70	93	15	53	22	83	140	151	119	144	107	39	149	143	146	103	82
##	[37]	106	76	6	55	7	105	85	33	84	80	153	101	87	54	36	10	111	91
##	[55]	68	108	58	61	98	122	79	74	97	1	147	116	77	152	113	37	159	24
##	[73]	63	96	9	88	40	14	62	154	100	71	49	3	136	27	47	56	50	73
##	[91]	72	19	31	81	67	137	35	21	69	13	16	95	160	135	156	32	112	28
##	[109]	8	89	90	86	110	157	75	57	102	126	134	138	65	48	66	141	45	11
##	[127]	125	78	5	2	155	52	104	132	139	25	109	43	38	44	42	29	46	20
##	[145]	94	23	34	41	12	26												

### 1.2 Data Cleaning

Independent variable 'maxsqfoot' is removed because

- We don't need two variables representing the size of the property. 'lotsize'is another area variable.
- We choose 'lotsize' instead of 'maxsqfoot' because there are 98 missing values for 'maxsqfoot'.

Also, we remove 11 data sets containing 'na' (missing values).

## II. Exploratory Data Analysis

### 2.1 Classify Variables

- Categorical Variable(s): location
- Discrete Variable(s): Number of bedrooms(bedroom), Number of bathrooms(bathroom), Number of parking spots(parking), Maximum square footage(maxsqfoot, removed)
- Continuous Variables: Sale price(sale), List price(list), Property tax(taxes), Lot size(lotsize), Frontage(lotwidth, removed), Length(lotlength, removed)

### 2.2 Correlation Matrix

TABLE 2.1: Correlation Coefficient Matrix

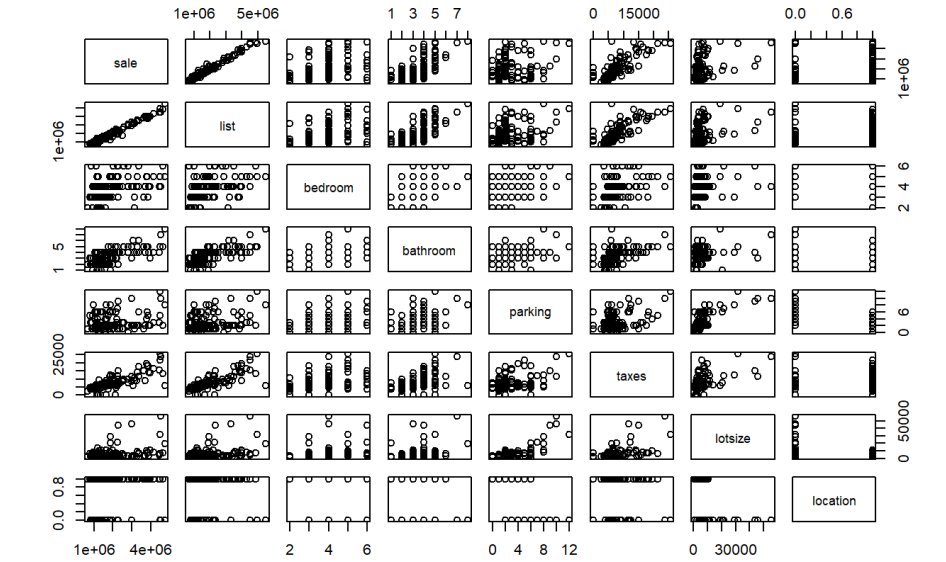
	sale	list	bedroom	bathroom	taxes	parking	lotsize	location
sale	1							
list	0.9874	1						
bedroom	0.4395	0.4436	1					
bathroom	0.6788	0.6955	0.539	1				
taxes	0.8087	0.7926	0.3931	0.5251	1			
parking	0.2486	0.2874	0.3543	0.3654	0.4287	1		
lotsize	0.3714	0.3873	0.2726	0.3175	0.5463	0.741	1	
location	0.0948	0.0585	-0.1491	-0.163	-0.1291	-0.771	-0.5589	1

Highest to Lowest correlation coefficient with 'sale price':

1. Last list price (list): 0.9874,
2. Property tax (taxes): 0.8087,
3. Number of bathrooms (bathroom): 0.6788,
4. Number of bedrooms (bedroom): 0.4395,
5. Lotsize (lotsize): 0.3714,
6. Number of parking spots (parking): 0.2486
7. Location of the property (location): 0.0948

### 2.2 Scatterplot Matrix

FIGURE 2.1 Scatterplot Matrix (1548)



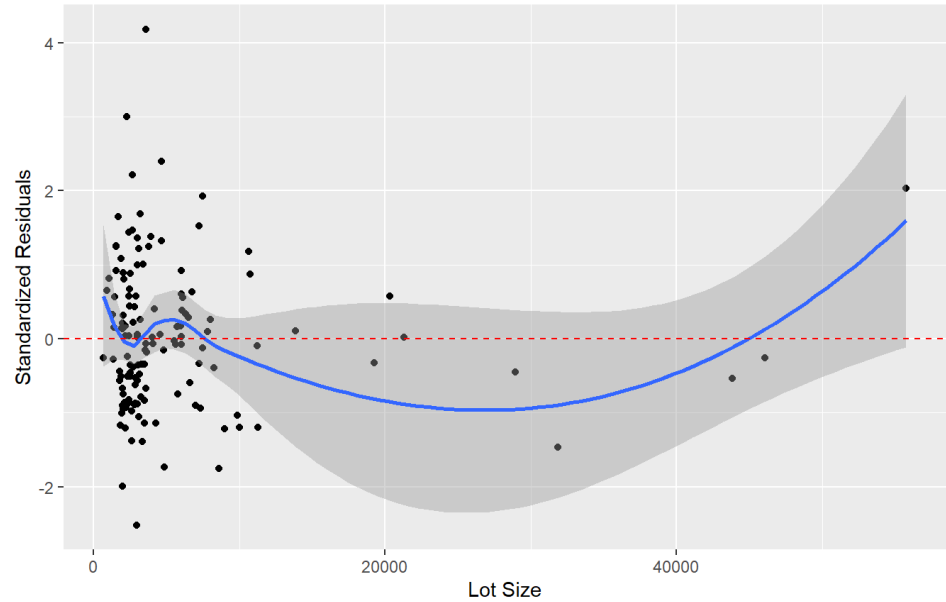
2.2.1 Violate The Constant Variance Assumption

Take a look at the first row except the dummy variable “location”; the scatter plots are approximately positively related except for the “lotsize”. Points on the last graph are centered at the left-bottom corner.

Thus, we guess that the data of ‘lotsize’ violate the constant variance assumption.

2.2.2 The Standardized Residuals Plot of ‘lotsize’

FIGURE 2.2 Standardized Residuals vs Lot Size (1548)



The residuals do not roughly form a horizontal band around the zero line (red), suggesting that the error terms’ variance are not constant. The conclusion proves that the constant variance assumption is not satisfied.

### III. Methods and Model

3.1 Fitted Linear Regression Model

3.1.1 Summary Table

TABLE 2.2: Summary of the Linear Regression Model  
Coefficient

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	100436.1181	68988.8890	1.4558	0.1478
list	0.8325	0.0242	34.3528	0.0000
bedroom	11935.7999	14318.8585	0.8336	0.4060
bathroom	7379.2088	14702.6122	0.5019	0.6166
parking	-21769.6909	10664.8601	-2.0413	0.0432
taxes	22.2828	4.9416	4.5093	0.0000
lotsize	1.5093	2.4429	0.6179	0.5377
location	58072.3904	52089.9799	1.1148	0.2670

3.1.2 Fitted Multiple Linear Regression Model (Full)

We get the full model given by:

$sale = 100436.12 + 0.83 \times list + 11935.80 \times bedroom + 7379.21 \times bathroom - 21769.70 \times parking + 22.28 \times taxes + 1.51 \times lotsize + 58072.40$

Note: location = 1 if the property is in Toronto, otherwise 0.

3.1.3 Significance of Variables

List price, parking spots, and property taxes are significant because the p-values of the three variables are smaller than the significance level 0.05. We are able to reject the null hypothesis that the coefficients of these variables are zero.

Keep all other variables constant, as the list price increased by 1, the property's average sale price will increase by 0.83. Keep all other variables constant, as the parking spots increased by 1, the property's average sale price will decrease by 21769.7. Keep all other variables constant, as the property tax increased by 1, the property's average sale price will increase by 22.28.

3.2 Find A Parsimonious Model (Backward, AIC)

We get the final model according to the AIC values using backward regression is given by

$$sale = 195400 + 0.85 \times list - 25200.70 \times parking + 22.6 \times taxes$$

3.3 Find A Parsimonious Model (Backward, BIC)

We get the final model according to the BIC values using backward regression is given by

$$sale = 195400 + 0.85 \times list - 25200.70 \times parking + 22.6 \times taxes$$

3.4 Check Multicollinearity

Global F-test for the final model is significant, equals to 2187. One of the partial F-test of three variable, list price, is much more significant than others. This indicates there does not exist multicollinearity, but one variable is more significant than other two.

3.5 Summary

The final model is given by:

$$sale = 195400 + 0.85 \times list - 25200.70 \times parking + 22.60 \times taxes$$

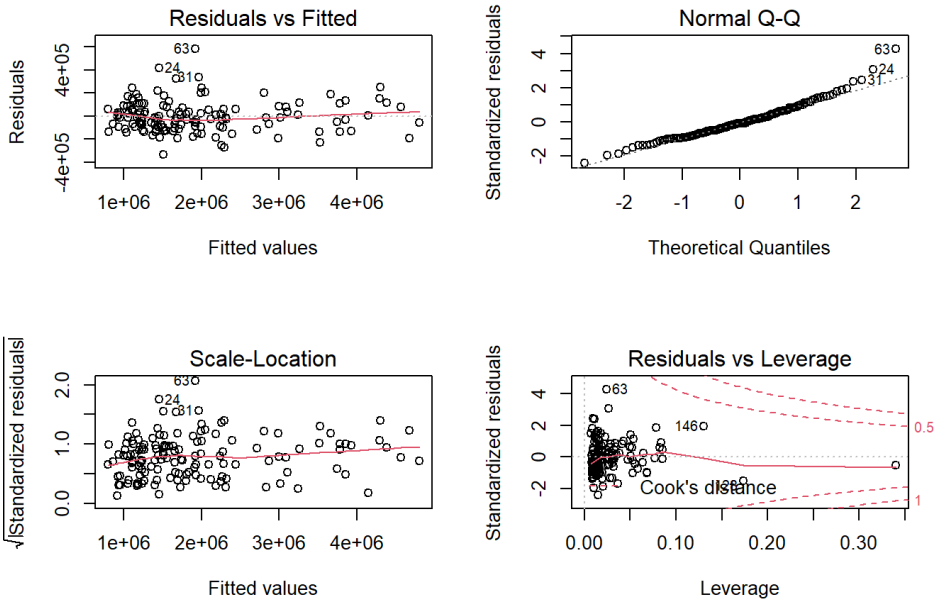
Keep all other variables constant, as the list price increased by 1, the property's average sale price will increase by 0.85. Keep all other variables constant, as the parking spots increased by 1, the property's average sale price will decrease by 25200.7. Keep all other variables constant, as the property tax increased by 1, the property's average sale price will increase by 22.6.

AIC or BIC produces the same model, but different from the full model. Three of the variables are removed, and the coefficients of the remaining three are slightly different. AIC and BIC are both penalized-likelihood criteria. The model given by AIC and BIC is more likely to approximate the true model. In this case, they choose three variables considered to be more influential to the dependent variable. The coefficients are different because three independent variables are deleted from the model.

IV. Discussions and Limitations

4.1 Diagnostic Plots

FIGURE 4.1 Diagnostic Plots (1548)



4.2 Interpretation Diagnostic Plots

4.2.1 Residual vs Fitted

In this case, the residuals are approximately randomly distributed around the horizontal zero line. This indicates that the residuals and the fitted values are uncorrelated. The assumption of equal variance (homoscedasticity) is satisfied.

4.2.2 Normal Q-Q

Most of the points follow the theoretical normal line. This indicates that the residuals are normally distributed. The assumption of the normal error MLR is satisfied.

4.2.3 Scale-Location

There is no obvious pattern shown in the graph. This indicates that the residuals are spread equally along with the ranges of predictors. The assumption of equal variance (homoscedasticity) is satisfied.

4.2.4 Residuals vs Leverage

All of the points are in the red line region. This indicates that there is no leverage point or outlier.

4.3 Future Work

- We could use box cox transformation to improve our model by making some of the variables more normally distributed.
- We could add more independent variables that may affect the sale price of properties.