

MA206 PD Meeting

Kevin Cummiskey

March 24, 2020

Review

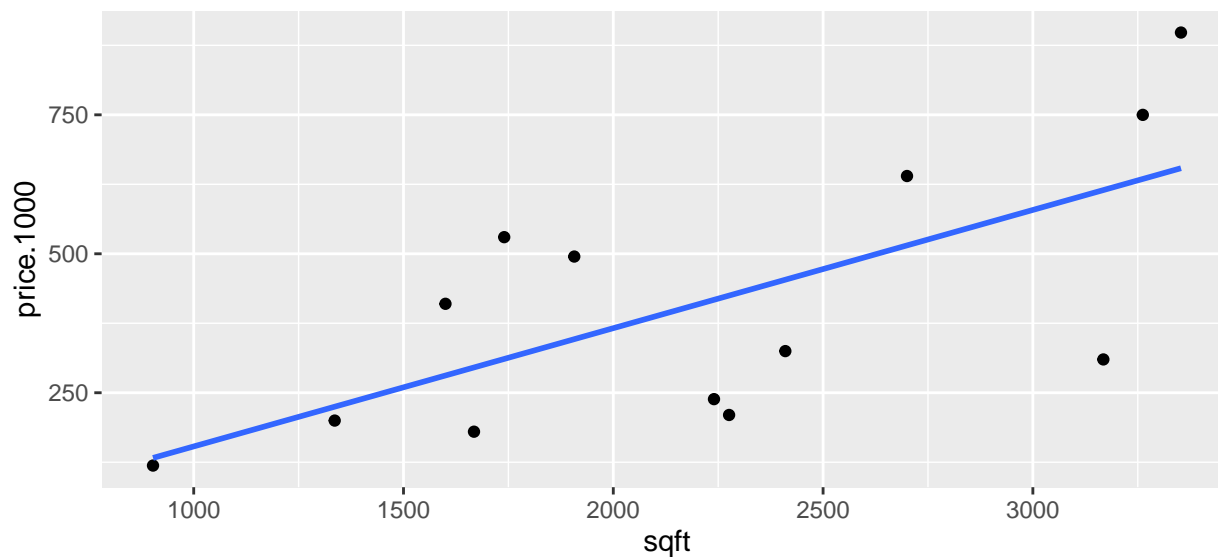
In Chapter 4, we've been looking at housing prices in Michigan. Specifically, we have been assessing the relationship between house size (sqft) and price (in \$100,000s).

$$\hat{price}_i = \beta_0 + \beta_1 sqft_i$$

where $price_i$ is the price of house i and $sqft_i$ is the size (sq ft) of house i .

```
houses = read.table(file = "http://www.isi-stats.com/isi2/data/housing.txt",
                    header = TRUE)
houses %>%
  ggplot(aes(x = sqft, y = price.1000)) +
  geom_point() +
  geom_smooth(method = "lm",
             se = FALSE)
```

```
## `geom_smooth()` using formula 'y ~ x'
```



```
model_simple <- lm(price.1000 ~ sqft, data = houses)
summary(model_simple)
```

```
##
## Call:
## lm(formula = price.1000 ~ sqft, data = houses)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -304.70 -128.44  -13.74   128.98   244.04
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -59.36870  161.36807  -0.368   0.7199
## sqft         0.21274    0.06963   3.055   0.0109 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 185.1 on 11 degrees of freedom
## Multiple R-squared:  0.4591, Adjusted R-squared:  0.4099
## F-statistic: 9.335 on 1 and 11 DF,  p-value: 0.01094
```

Based on this model, estimate the price per square foot.

Briefly explain why location (lake or nonlake) may be a confounding variable.

We can adjust for a confounding variable by including it in a multiple regression model. In this model, we estimate the price per square foot when we hold the location of the house constant.

$$\hat{price}_i = \alpha_0 + \alpha_1 sqft_i + \alpha_2 lake_i$$

where $price_i$ is the price of house i , $sqft_i$ is the size (sq ft) of house i , and $lake_i$ is 1 if house i is lakefront and is 0 otherwise.

What assumptions does this model make?

How do we interpret α_1 ? α_2 ?

In general, when will α_1 from this model equal β_1 in the simple model?

```
#Reverse coding of lake so 1 is lake and 0 is not lake
houses$lake = factor(houses$lake, levels = c("notlakefront",
                                             "lakefront"))
```

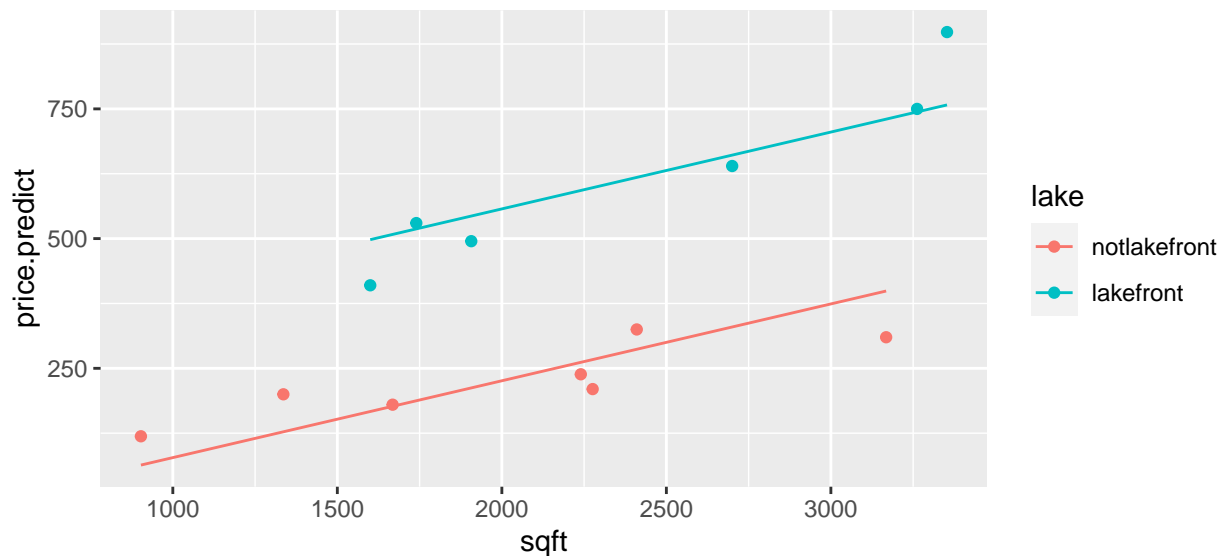
```

model_withLake = lm(price.1000 ~ sqft + lake, data = houses)

houses %>%
  mutate(price.predict = predict(model_withLake, newdata = .),
         residuals = price.1000 - price.predict) -> houses

houses %>%
  ggplot(aes(x = sqft, y = price.predict, color = lake)) +
  geom_line() +
  geom_point(aes(y = price.1000))

```



```
summary(model_withLake)
```

```

##
## Call:
## lm(formula = price.1000 ~ sqft + lake, data = houses)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -89.059 -48.444   3.072  38.191 140.421
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -70.1821    62.8062  -1.117  0.289933
## sqft           0.1481     0.0283   5.233 0.000383 ***
## lakelakefront 331.2235    41.8470   7.915 1.29e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 72.02 on 10 degrees of freedom
## Multiple R-squared:  0.9255, Adjusted R-squared:  0.9106
## F-statistic: 62.15 on 2 and 10 DF, p-value: 2.289e-06

```

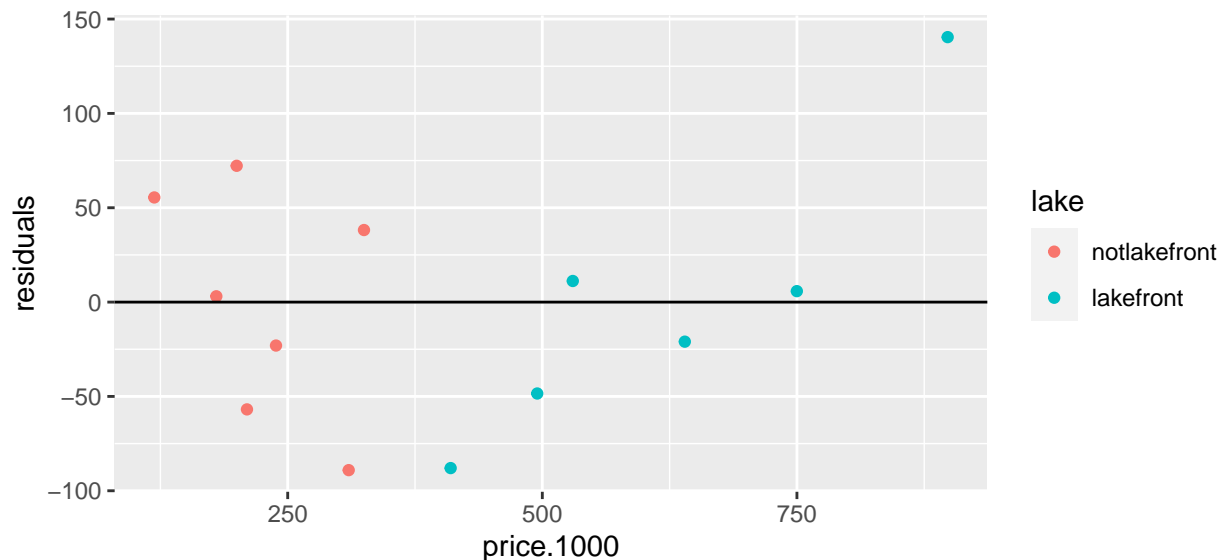
How did the estimate of the price per square footage change when we held location constant?

Calculate the expected price of a 2500 sq ft house that is not on the lake front.

Calculate the expected price of a 2500 sq ft house on the lake front.

Let's take a look at the residuals vs the predicted (fitted) values.

```
houses %>% ggplot(aes(x = price.1000,  
                      y = residuals,  
                      color = lake)) +  
  geom_point() + geom_hline(yintercept = 0)
```



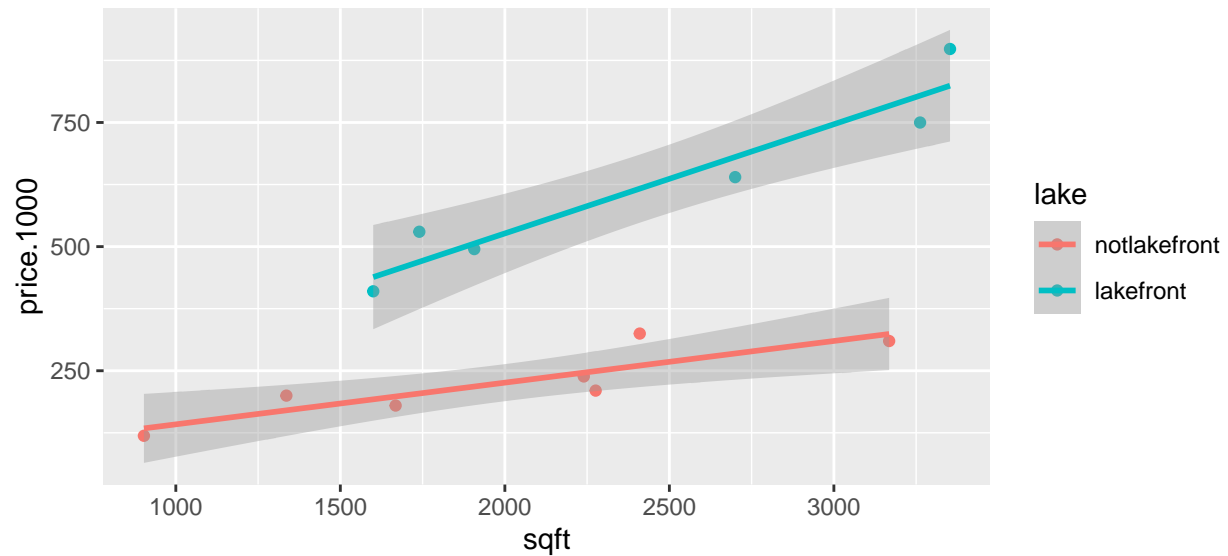
Do we see evidence of an interaction here?

Two models

If we suspect the price per square foot differs for lake front and nonlake front, we could fit two separate simple models.

```
houses %>% ggplot(aes(x = sqft, y = price.1000, color = lake)) +  
  geom_point() + geom_smooth(method = "lm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```



```
lake_model <- lm(price.1000 ~ sqft, data = houses %>% filter(lake == "lakefront"))
coef(lake_model)
```

```
## (Intercept)      sqft
## 86.7643819    0.2198952
```

```
nonlake_model <- lm(price.1000 ~ sqft, data = houses %>% filter(lake == "notlakefront"))
coef(nonlake_model)
```

```
## (Intercept)      sqft
## 58.11340672    0.08394444
```

Interactions

Let's look at a model with an interaction.

$$\hat{price}_i = \beta_0 + \beta_1 sqft_i + \beta_2 lake_i + \beta_3 sqft_i lake_i$$

How do we interpret β_1 ? β_2 ? β_3 ? $\beta_1 + \beta_3$?

```
model_interaction = lm(price.1000 ~ sqft * lake, data = houses)
summary(model_interaction)
```

```
##
## Call:
## lm(formula = price.1000 ~ sqft * lake, data = houses)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -54.16  -28.60  -14.15   29.64   73.93
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    58.11341    56.68270     1.025  0.33202
## sqft           0.08394     0.02675     3.138  0.01197 *
## lakelakefront  28.65098    91.32560     0.314  0.76088
## sqft:lakelakefront 0.13595     0.03895     3.491  0.00682 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 49.48 on 9 degrees of freedom
## Multiple R-squared:  0.9684, Adjusted R-squared:  0.9578
## F-statistic: 91.84 on 3 and 9 DF, p-value: 4.547e-07
```

```
anova(model_interaction)
```

```
## Analysis of Variance Table
##
## Response: price.1000
##           Df Sum Sq Mean Sq F value    Pr(>F)
## sqft       1 319753  319753 130.611 1.166e-06 ***
## lake       1 324911  324911 132.718 1.089e-06 ***
## sqft:lake   1  29829   29829  12.184  0.006824 **
## Residuals   9  22033    2448
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Write the equation for the price of a lake front house as a function of size.

Write the equation for the price of a nonlake house as a function of size.

How do these equations compare to the two models above? Why is the interaction model preferable?

Is there evidence the price per square foot is different for lakefront and nonlakefront homes? How do you know?

What would you conclude from these results? Your answer should include discussion of effect sizes, significance, and overall predictive capability of the model.

The main effect of location is not significant in the model with an interaction. Should we conclude the location of the house is not associated with price? Justify your answer.