

Class work

Fentaw Abitew

2022-10-31

Researchers were interested in predicting residential home sales prices in a Midwestern city as a function of various characteristics of the home and surrounding property. Data on 522 transactions were obtained for home sales during the year 2002. The 13 variables are

- **Price:** Sales price of residence (in dollars)
- **Area:** Finished area of residence (in square feet)
- **Bed:** Total number of bedrooms in residence
- **Bath:** Total number of bathrooms in residence
- **AC:** 1 = presence of air conditioning, 0 = absence of air conditioning
- **Garage:** Number of cars that a garage will hold
- **Pool:** 1 = presence of a pool, 0 = absence of a pool
- **Year:** Year property was originally constructed
- **Quality:** Index for quality of construction. **High**, **Medium**, or **Low**.
- **Style:** Categorical variable indicating architectural style
- **Lot:** Lot size (in square feet)
- **Highway:** 1 = highway adjacent, 0 = highway not adjacent.

The data are available in “estate.csv” at https://dcgerard.github.io/stat_412_612/data/estate.csv.

Perform an exploratory data analysis to come up with some hypotheses. Some suggested ways to focus your research:

- Which variables are categorical? Which are quantitative?
- Change the values of the categorical variables to something more informative.
- What variables are marginally associated with price? Use plots and summary statistics.
- What variables are marginally associated with each other? Use plots and summary statistics.
- If a variable is marginally associated with price, are there some other variables that could explain that association? Use plots and summary statistics.
- Does there appear to be any discrete groupings of houses?
- Are there any unusual observations?
- What transformations should you perform to make associations more linear?
- Try making new variables based on existing variables.
- What variables should be discretized (or have values aggregated) because there are too few values and/or the association seems discrete?
- If you know linear regression, try out the `lm()` and `step()` functions to choose a tentative model.

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.2 --
## v ggplot2 3.3.6      v purrr   0.3.4
## v tibble  3.1.8      v dplyr  1.0.10
## v tidyr   1.2.1      v stringr 1.4.1
## v readr   2.1.2      v forcats 0.5.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
```

```
## x dplyr::lag()      masks stats::lag()
library(GGally) ## for pairs plot

## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2
theme_set(theme_bw())

estate<-read_csv('https://dcgerard.github.io/stat_412_612/data/estate.csv')

## Rows: 522 Columns: 12
## -- Column specification -----
## Delimiter: ","
## chr  (1): Quality
## dbl (11): Price, Area, Bed, Bath, AC, Garage, Pool, Year, Style, Lot, Highway
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
head(estate)

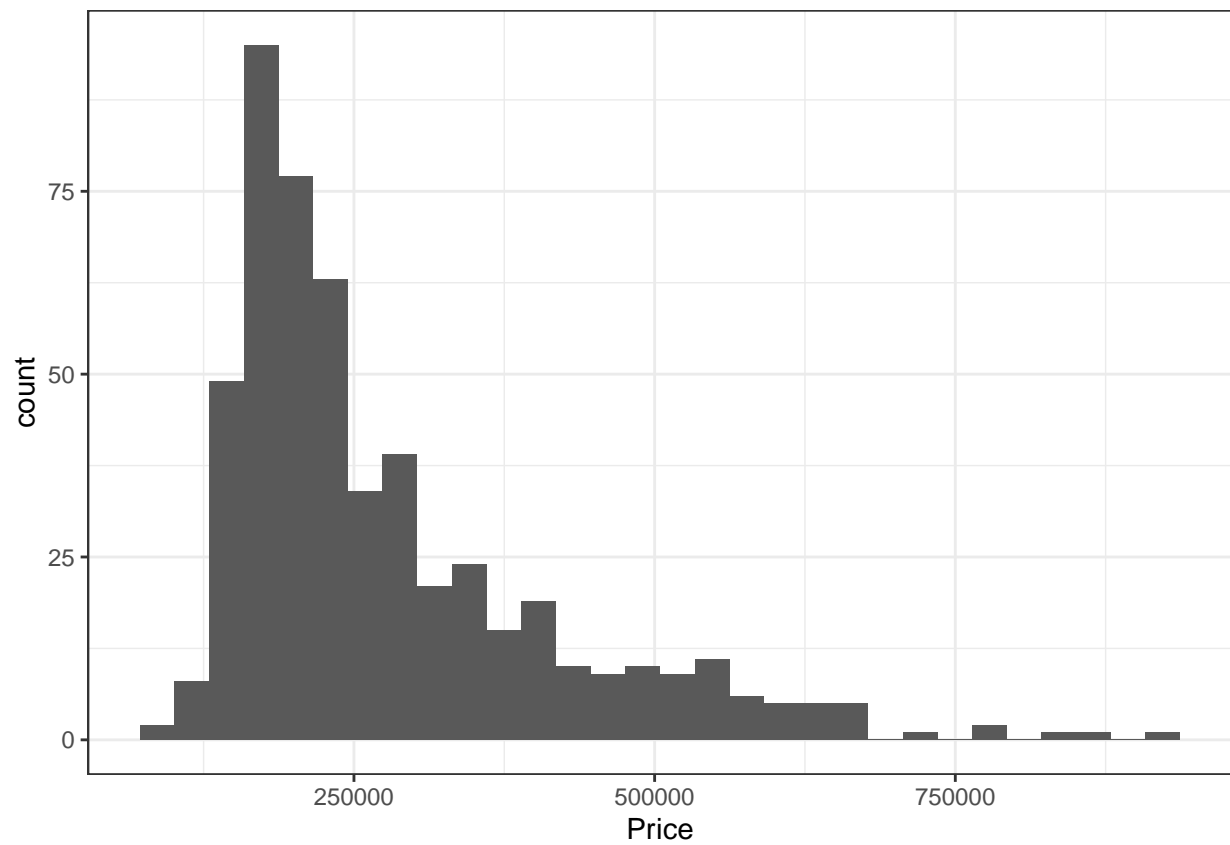
## # A tibble: 6 x 12
##   Price Area  Bed Bath  AC Garage Pool Year Quality Style  Lot Highway
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <chr>  <dbl> <dbl> <dbl>
## 1 360000 3032   4    4    1     2    0 1972 Medium    1 22221    0
## 2 340000 2058   4    2    1     2    0 1976 Medium    1 22912    0
## 3 250000 1780   4    3    1     2    0 1980 Medium    1 21345    0
## 4 205500 1638   4    2    1     2    0 1963 Medium    1 17342    0
## 5 275500 2196   4    3    1     2    0 1968 Medium    7 21786    0
## 6 248000 1966   4    3    1     5    1 1972 Medium    1 18902    0
```

I'm going to change the obvious variables that should be factors to factors.

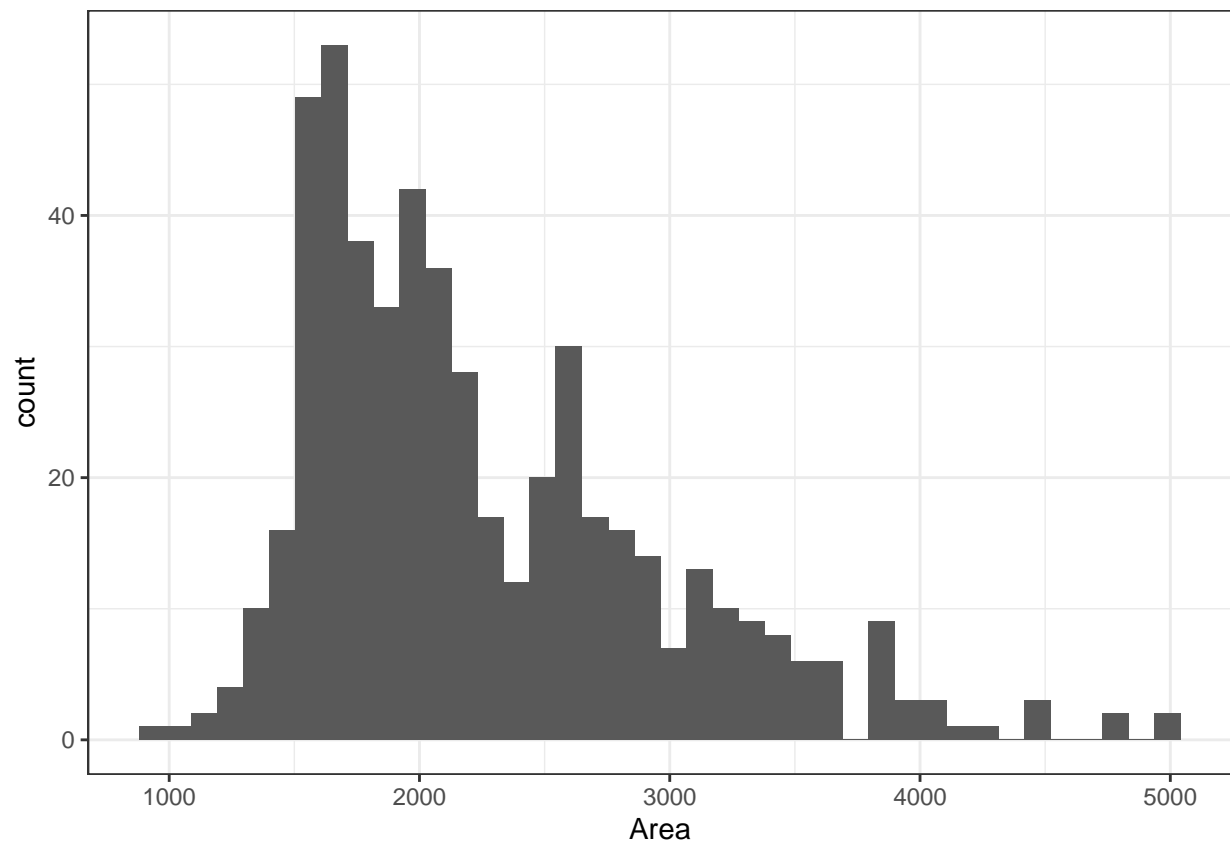
```
estate %>%
  mutate(AC = recode(AC, "1" = "AC", "0" = "noAC"),
         Pool = recode(Pool, "1" = "Pool", "0" = "noPool"),
         Style = as.factor(Style),
         Highway = recode(Highway, "1" = "Highway", "0" = "noHighway")) ->
estate
```

Plotting

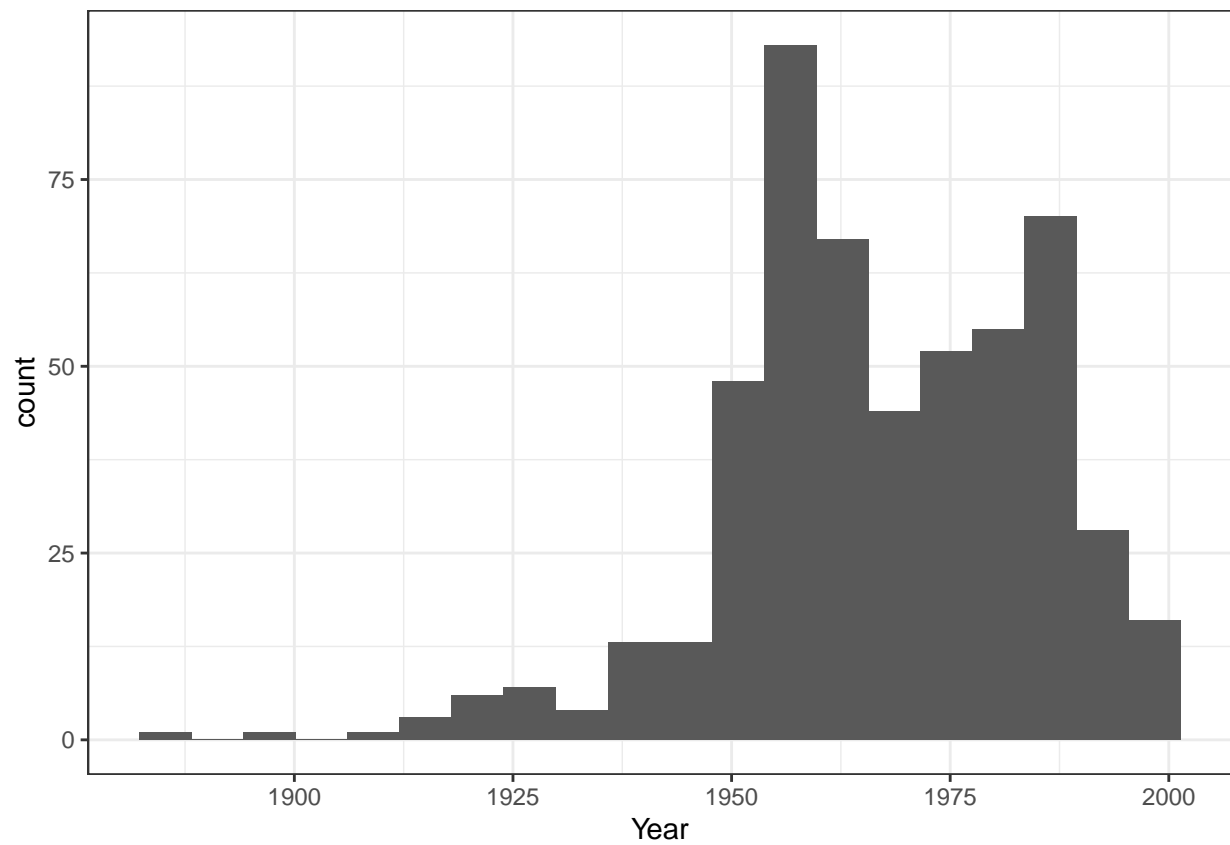
```
ggplot(estate, aes(x = Price)) +
  geom_histogram(bins = 30)
```



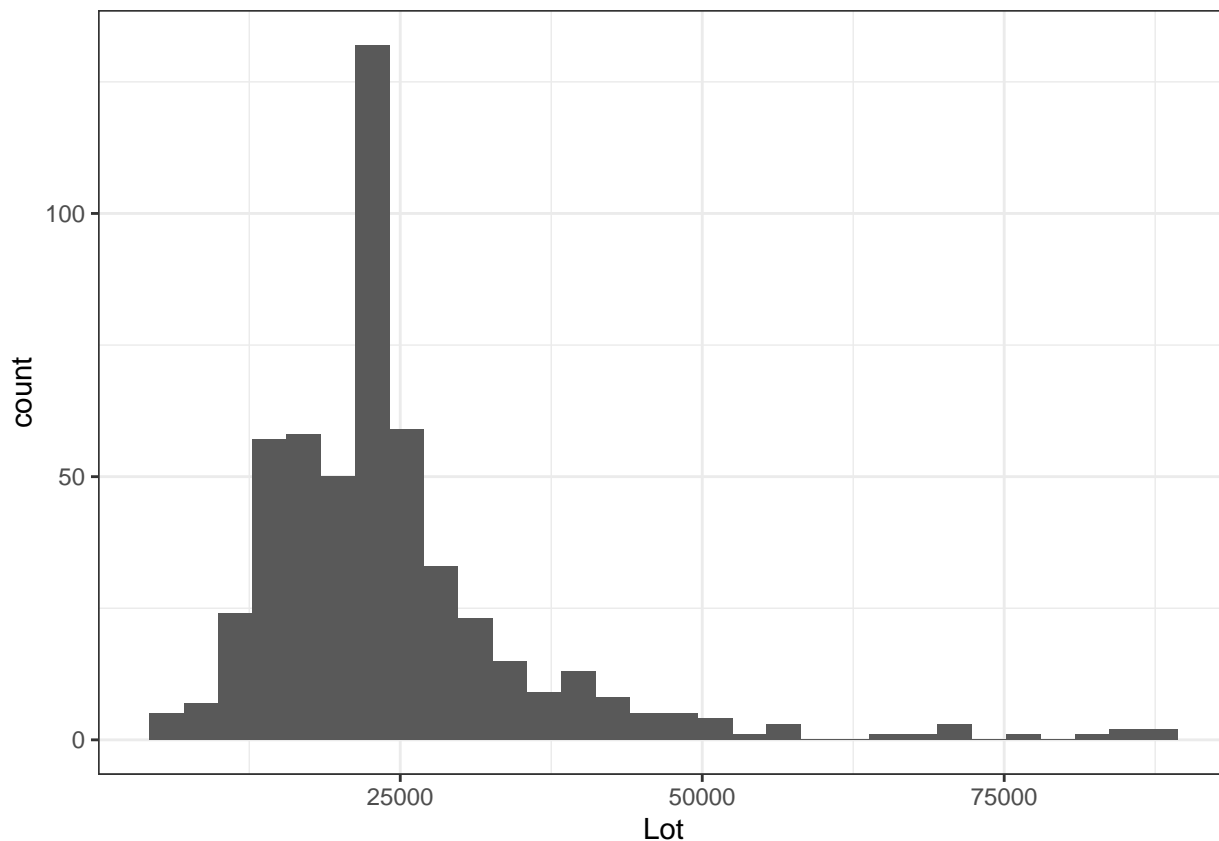
```
## Two bumps in Area followed by a long tail
ggplot(estate, aes(x = Area)) +
  geom_histogram(bins = 40)
```



```
ggplot(estate, aes(x = Year)) +  
  geom_histogram(bins = 20)
```



```
ggplot(estate, aes(x = Lot)) +  
  geom_histogram(bins = 30)
```



```
## mostly 3 and 4 bedroom houses, but there is a 0 bedroom house
## Is that a studio?
table(estate$Bed)
```

```
##
##  0  1  2  3  4  5  6  7
##  1  9 64 202 179 52 12  3
```

```
## A 0 bathroom house??? Is that the same house?
table(estate$Bath)
```

```
##
##  0  1  2  3  4  5  6  7
##  1 71 171 175 84 17  1  2
```

```
## Let's look at that unit
estate %>%
  filter(Bath == 0)
```

```
## # A tibble: 1 x 12
##   Price Area Bed Bath AC Garage Pool Year Quality Style Lot Highway
##   <dbl> <dbl> <dbl> <dbl> <chr> <dbl> <chr> <dbl> <chr> <fct> <dbl> <chr>
## 1 528750 2129 0 0 AC 3 noPool 1992 High 1 37414 noHigh~
```

```
## It's price is on the high end for having no bathroom!
## (92nd percentile)
```

```
estate %>%
  filter(Bath == 0) %>%
  select(Price) %>%
  c() ->
```

```

    weird_house_price
mean(estate$Price < weird_house_price)

## [1] 0.9214559
## I would keep in mind removing that house if I was to go on and do a
## Linear regression

table(estate$AC)

##
##   AC noAC
##  434   88
## One garage holds 7 cars?
table(estate$Garage)

##
##   0   1   2   3   4   5   7
##   7  52 353 106   2   1   1

table(estate$Pool)

##
## noPool   Pool
##   486     36

table(estate$Quality)

##
##   High   Low Medium
##    68    164    290

table(estate$Style)

##
##   1   2   3   4   5   6   7   9  10  11
## 214  58  64  11  18  18 136   1   1   1
## Very few houses on a highway. I would be
## Careful about inferences there
table(estate$Highway)

##
##   Highway noHighway
##        11        511

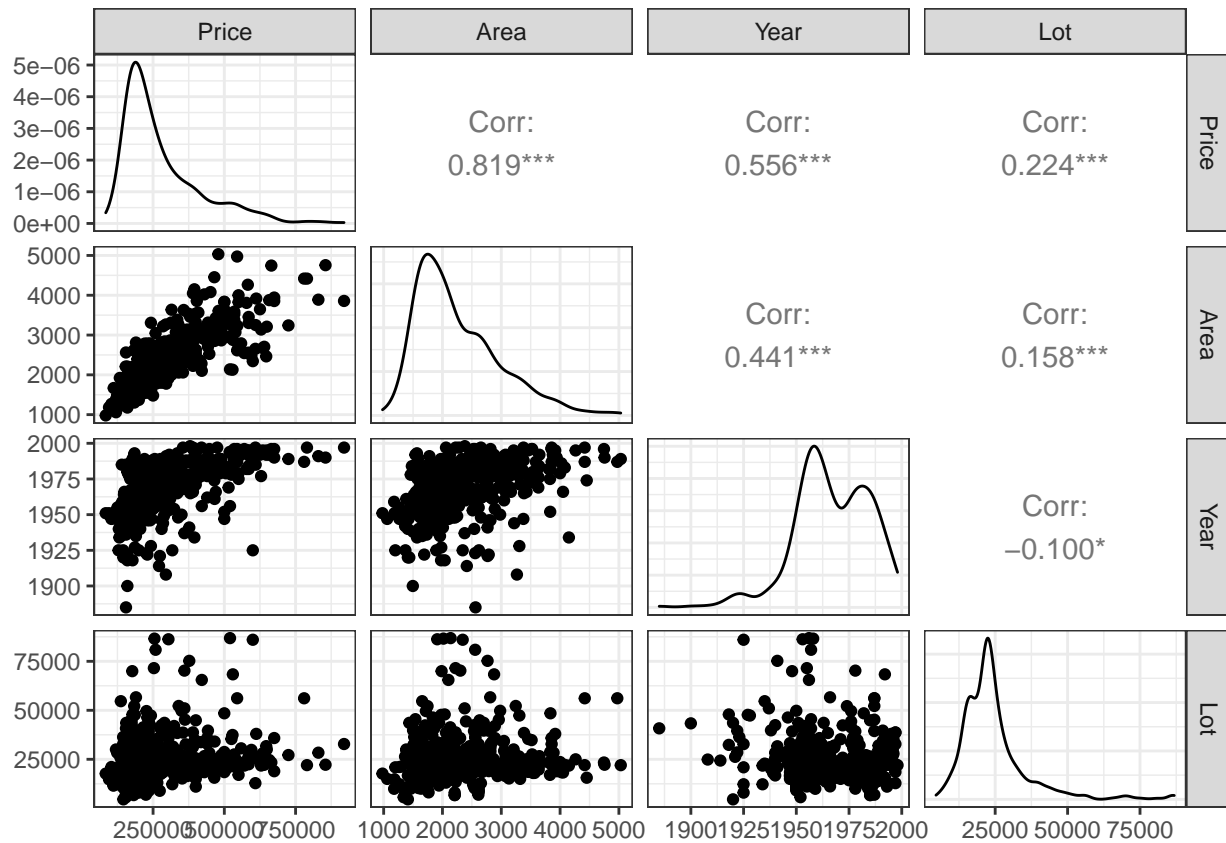
```

Look at bivariate associations

```

estate %>%
  select(Price, Area, Year, Lot) %>%
  ggpairs()

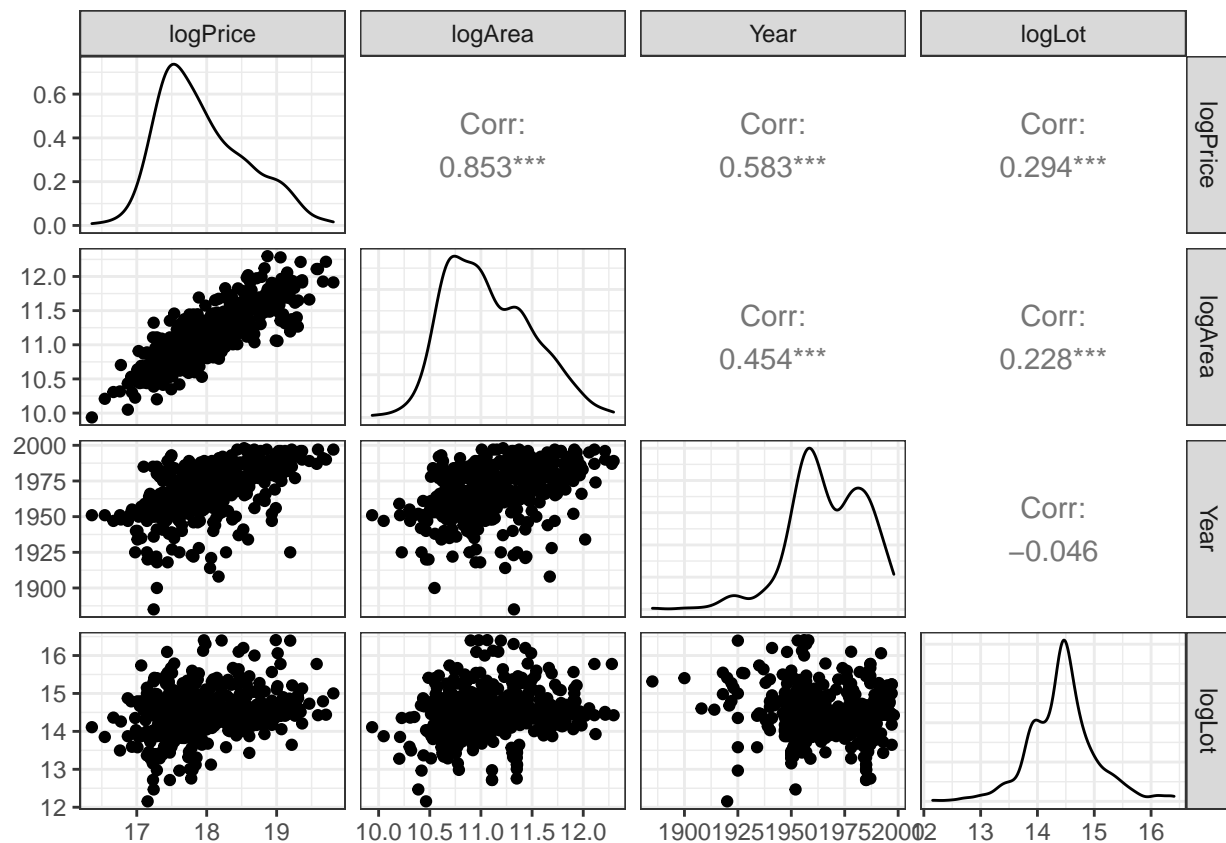
```



It seems that we can log price and area and lot pretty safely.

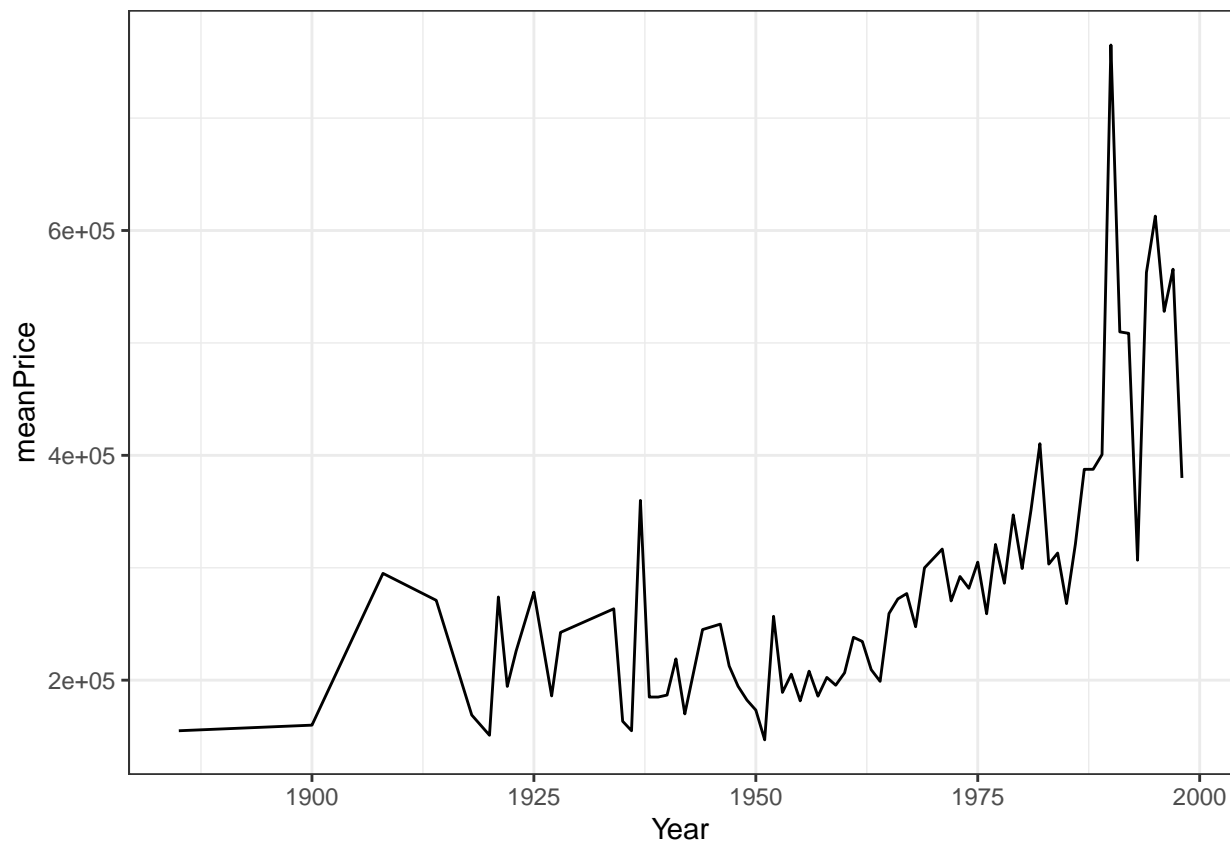
```
estate %>%
  mutate(logPrice = log2(Price), logArea = log2(Area), logLot = log2(Lot)) ->
  estate

estate %>%
  select(logPrice, logArea, Year, logLot) %>%
  ggpairs()
```

It seems that area has the strongest relationship to price

```
estate %>%
  group_by(Year) %>%
  summarize(meanPrice = mean(Price)) %>%
  ggplot(aes(x = Year, y = meanPrice)) +
  geom_line()
```



```
summary(estate$Year)
```

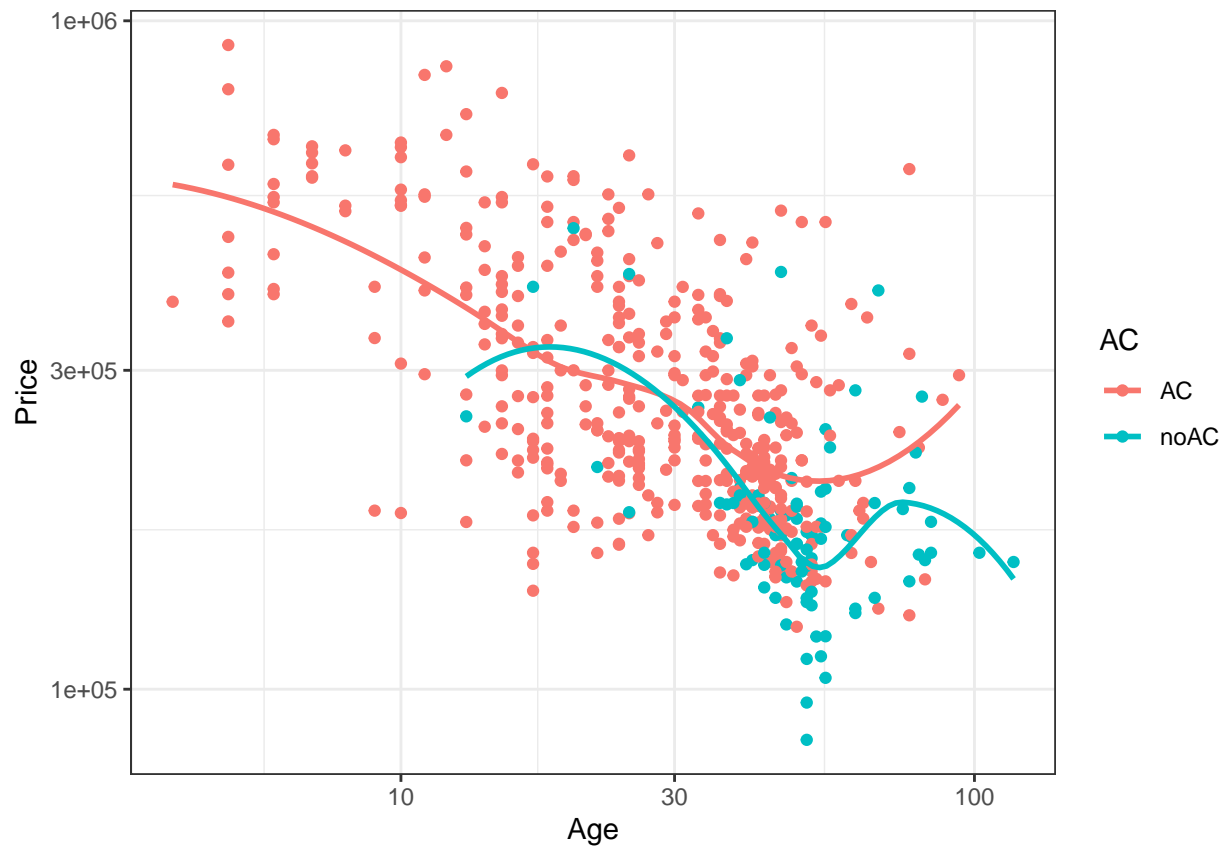
```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      1885   1956   1966   1967   1981   1998
```

```
## Define Age of hours at sale
estate %>%
  mutate(Age = 2002 - Year) ->
  estate
```

Q: Does anything variable get more important with the age of the house?

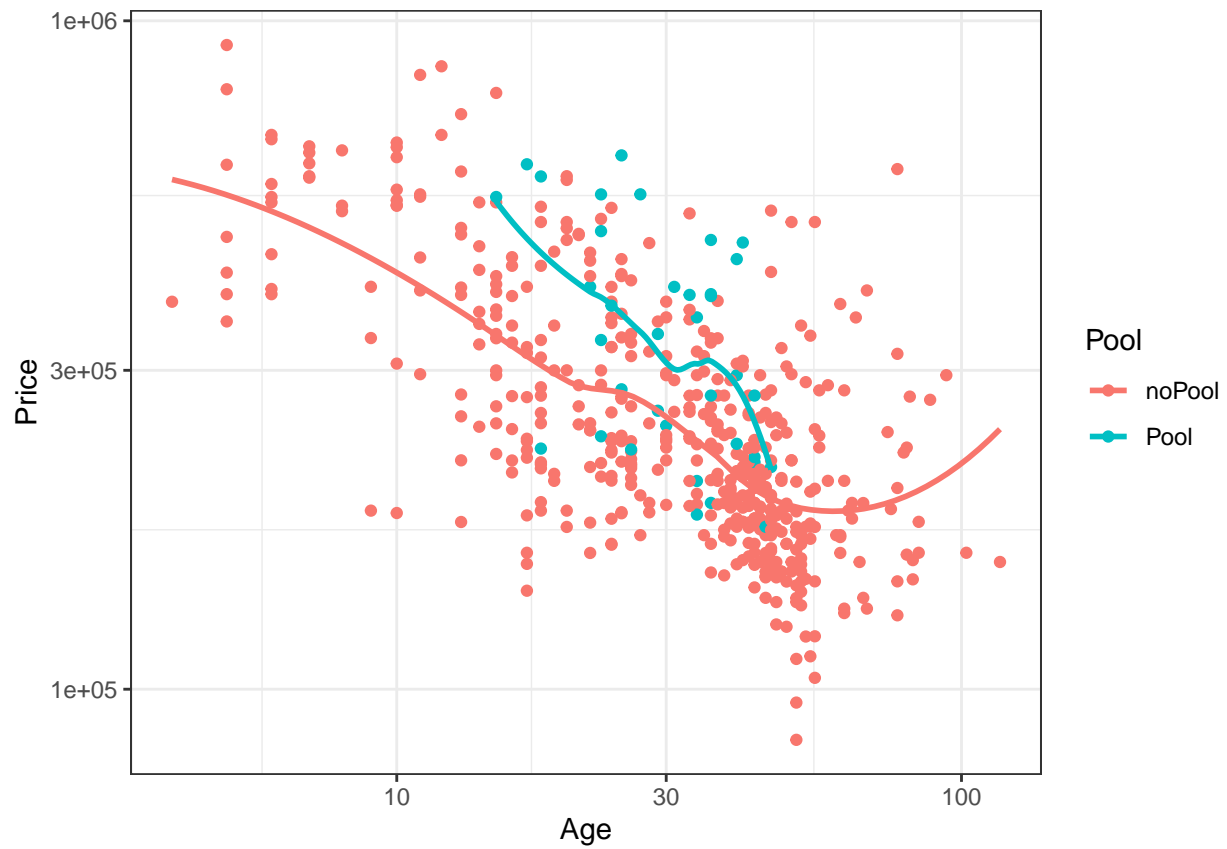
```
## Most no-ac houses are older. And once you adjust for age.
## But it still seems that there is an AC effect, which is particularly
## strong for older houses.
ggplot(estate, aes(x = Age, y = Price, color = AC)) +
  geom_point() +
  geom_smooth(se = FALSE) +
  scale_x_log10() +
  scale_y_log10()
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

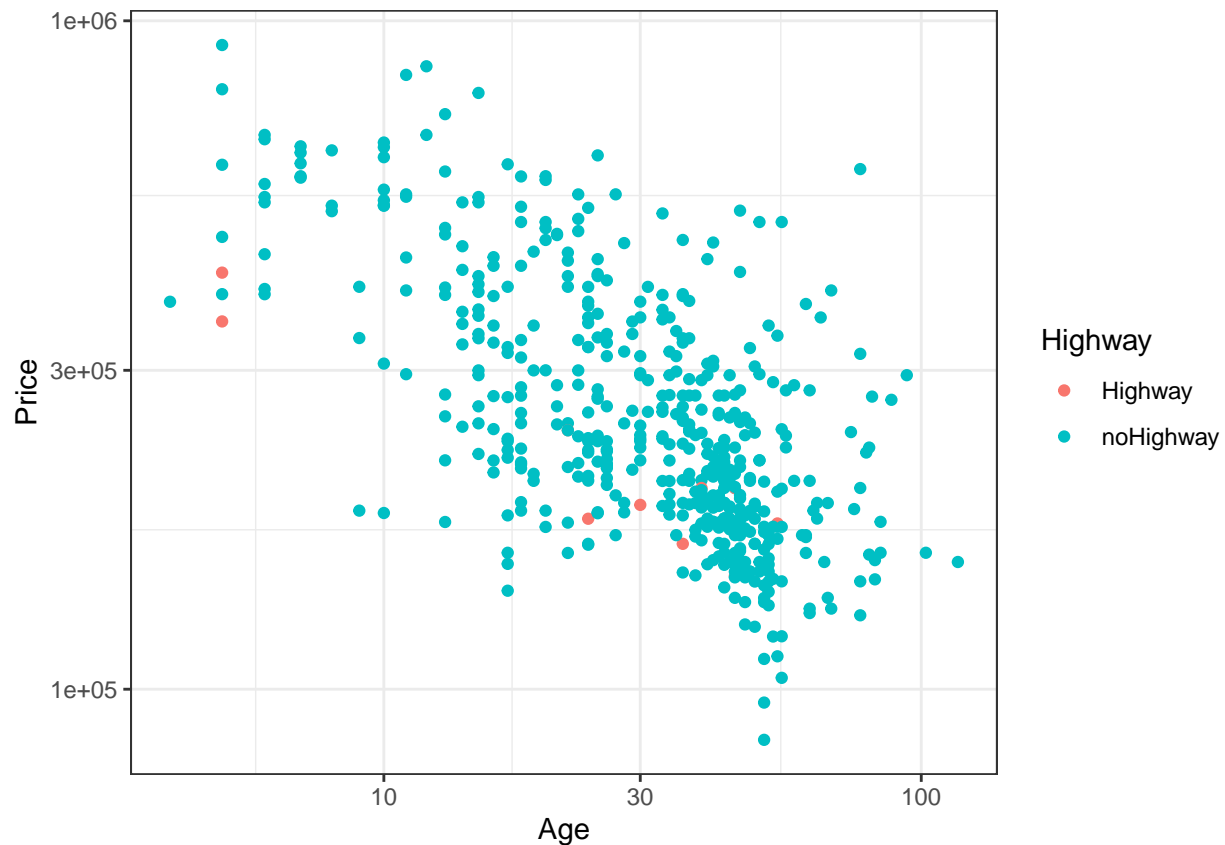


```
## There seems to be an additive pool effect
ggplot(estate, aes(x = Age, y = Price, color = Pool)) +
  geom_point() +
  geom_smooth(se = FALSE) +
  scale_x_log10() +
  scale_y_log10()

## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```



```
## The houses near the highway almost always have a lower price
ggplot(estate, aes(x = Age, y = Price, color = Highway)) +
  geom_point() +
  scale_x_log10() +
  scale_y_log10()
```

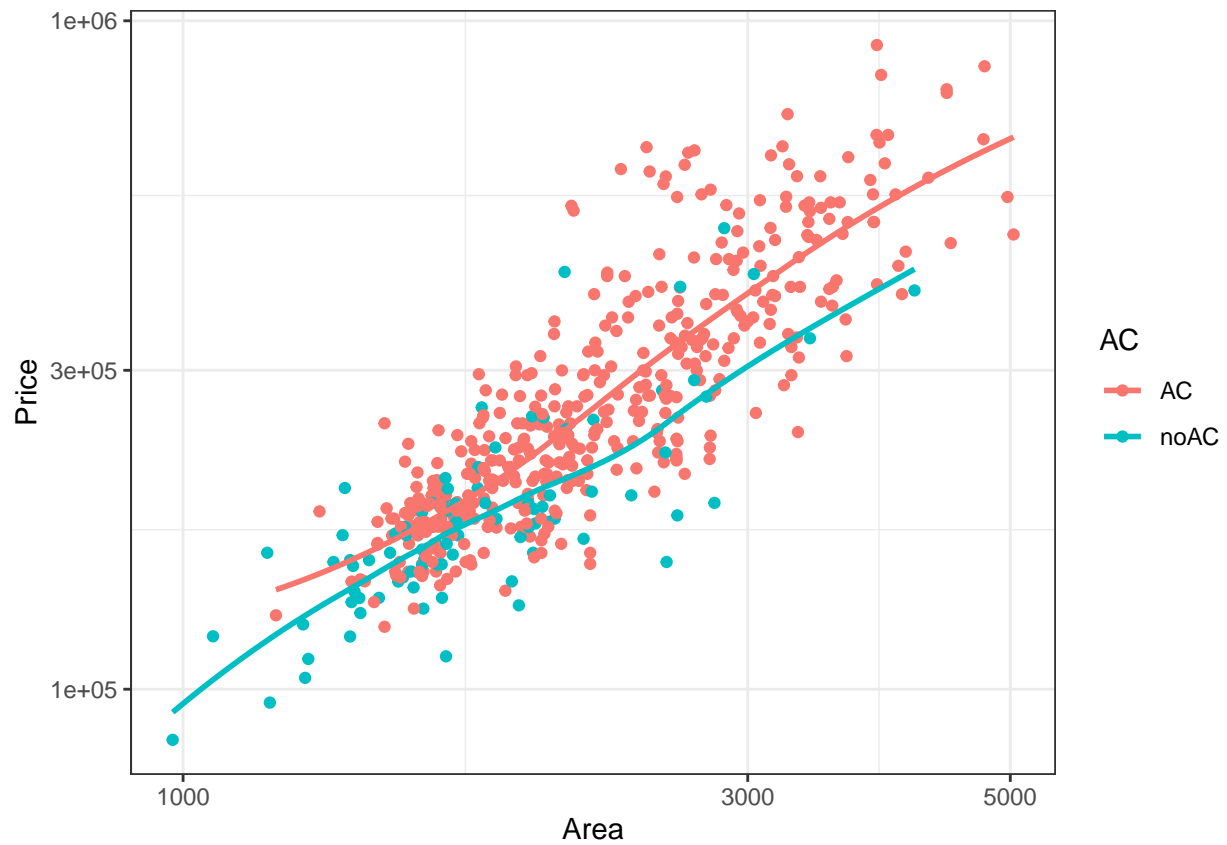


```
estate %>%
  mutate(logAge = log2(Age)) ->
  estate
```

Redo coloring with Area

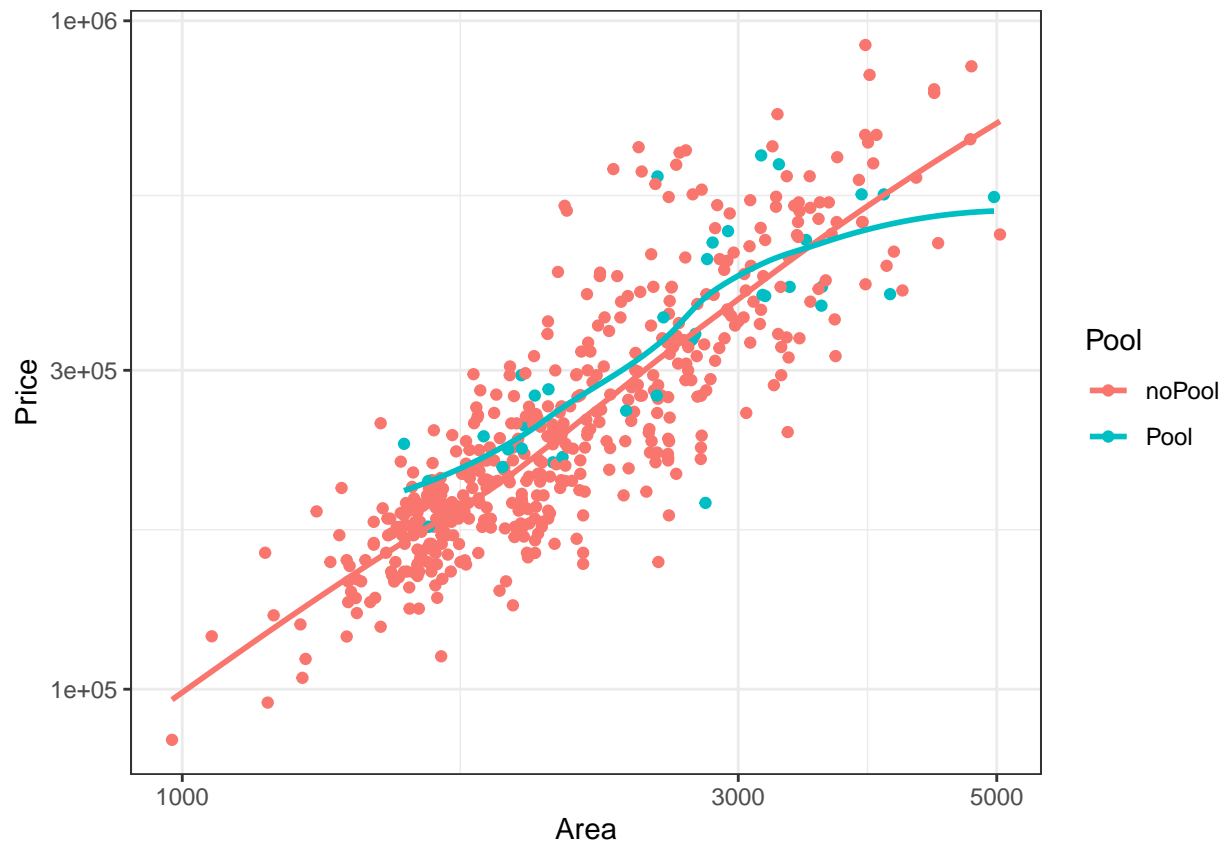
```
## Most no-ac houses are older. And once you adjust for age.
## But it still seems that there is an AC effect, which is particularly
## strong for older houses.
ggplot(estate, aes(x = Area, y = Price, color = AC)) +
  geom_point() +
  geom_smooth(se = FALSE) +
  scale_x_log10() +
  scale_y_log10()

## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

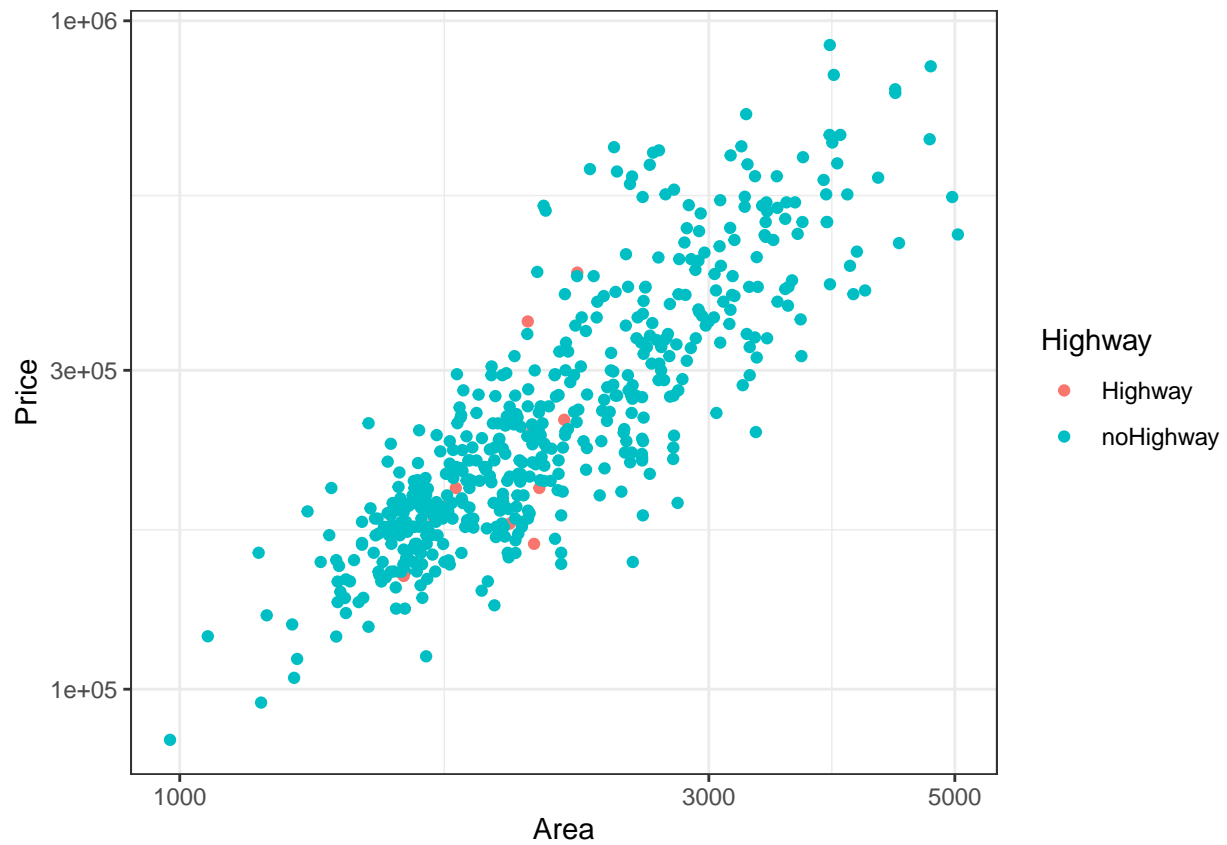


```
## There seems to be an additive pool effect
ggplot(estate, aes(x = Area, y = Price, color = Pool)) +
  geom_point() +
  geom_smooth(se = FALSE) +
  scale_x_log10() +
  scale_y_log10()

## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

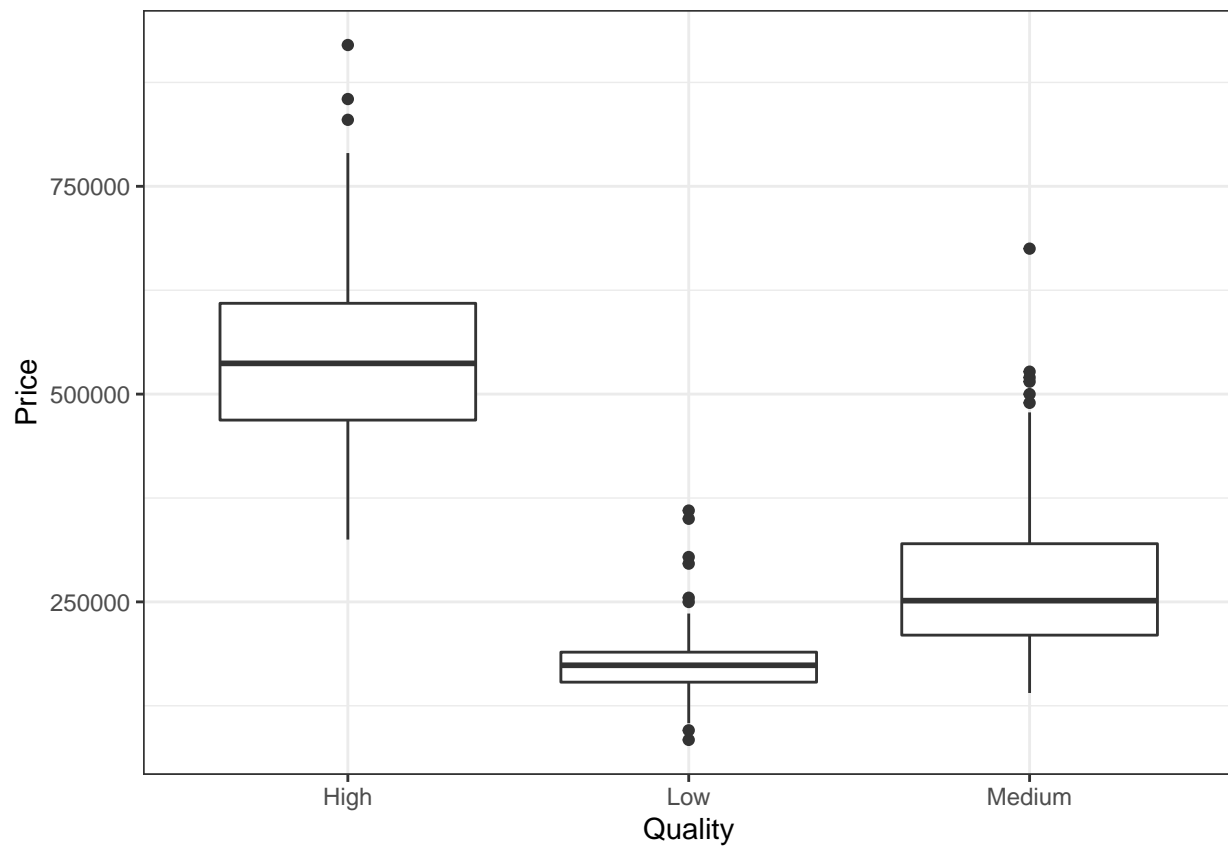


```
## The houses near the highway almost always have a lower price
ggplot(estate, aes(x = Area, y = Price, color = Highway)) +
  geom_point() +
  scale_x_log10() +
  scale_y_log10()
```

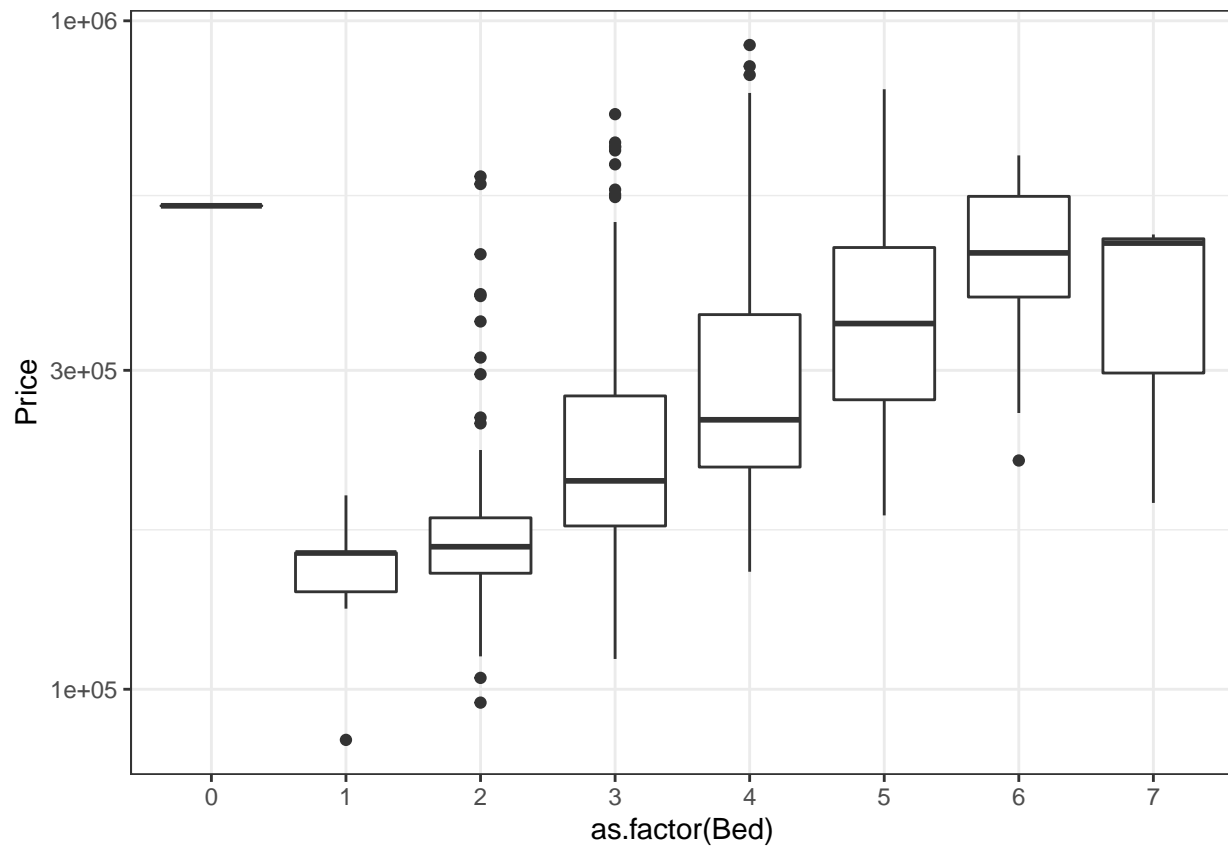


Look at the price by those other categories

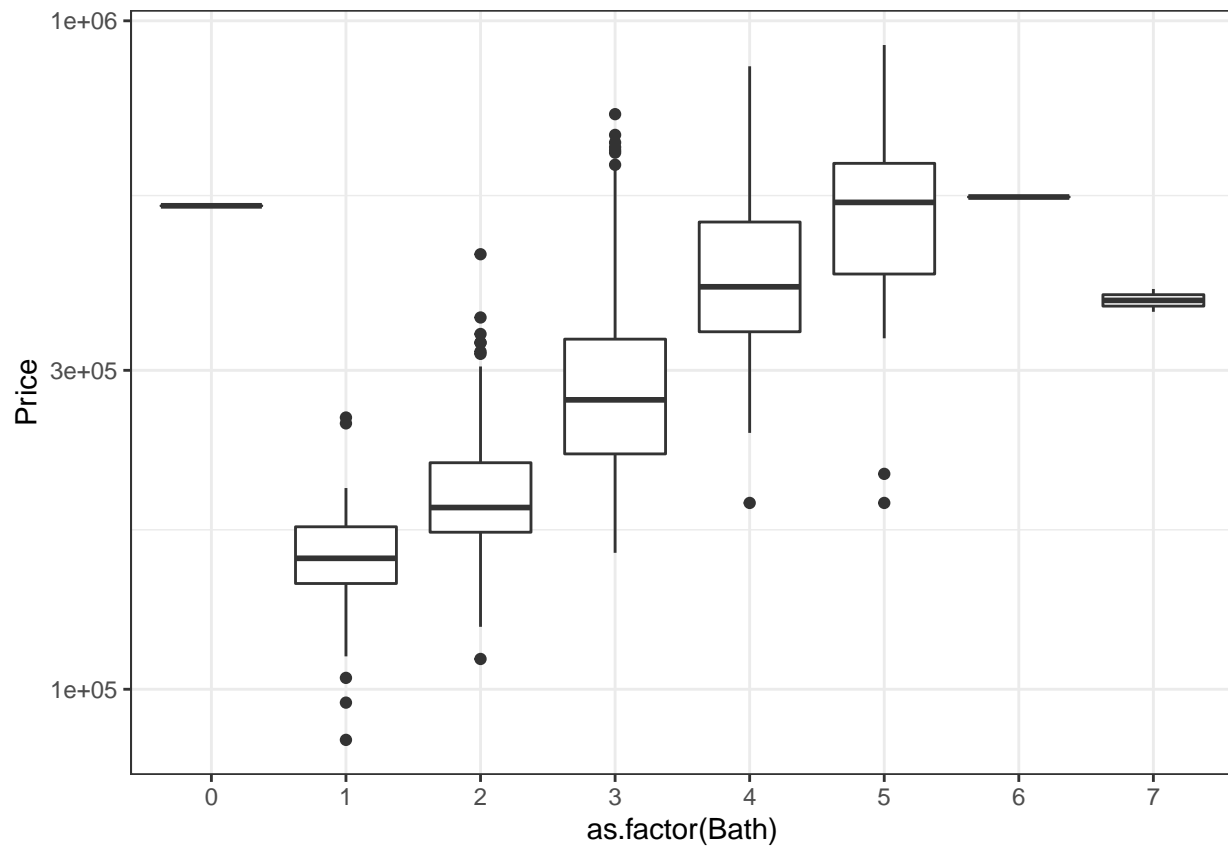
```
estate %>%  
  ggplot(aes(x = Quality, y = Price)) +  
  geom_boxplot()
```

```
## Seems that we could treat bed as a  
## Quantitative variable if we got  
## rid of that 0 house  
estate %>%  
  ggplot(aes(x = as.factor(Bed), y = Price)) +  
  geom_boxplot() +  
  scale_y_log10()
```



```
## I'm going to marge 5 bath into 5 and above,
## But still treat it as a quantitative variable
estate %>%
  ggplot(aes(x = as.factor(Bath), y = Price)) +
  geom_boxplot() +
  scale_y_log10()
```

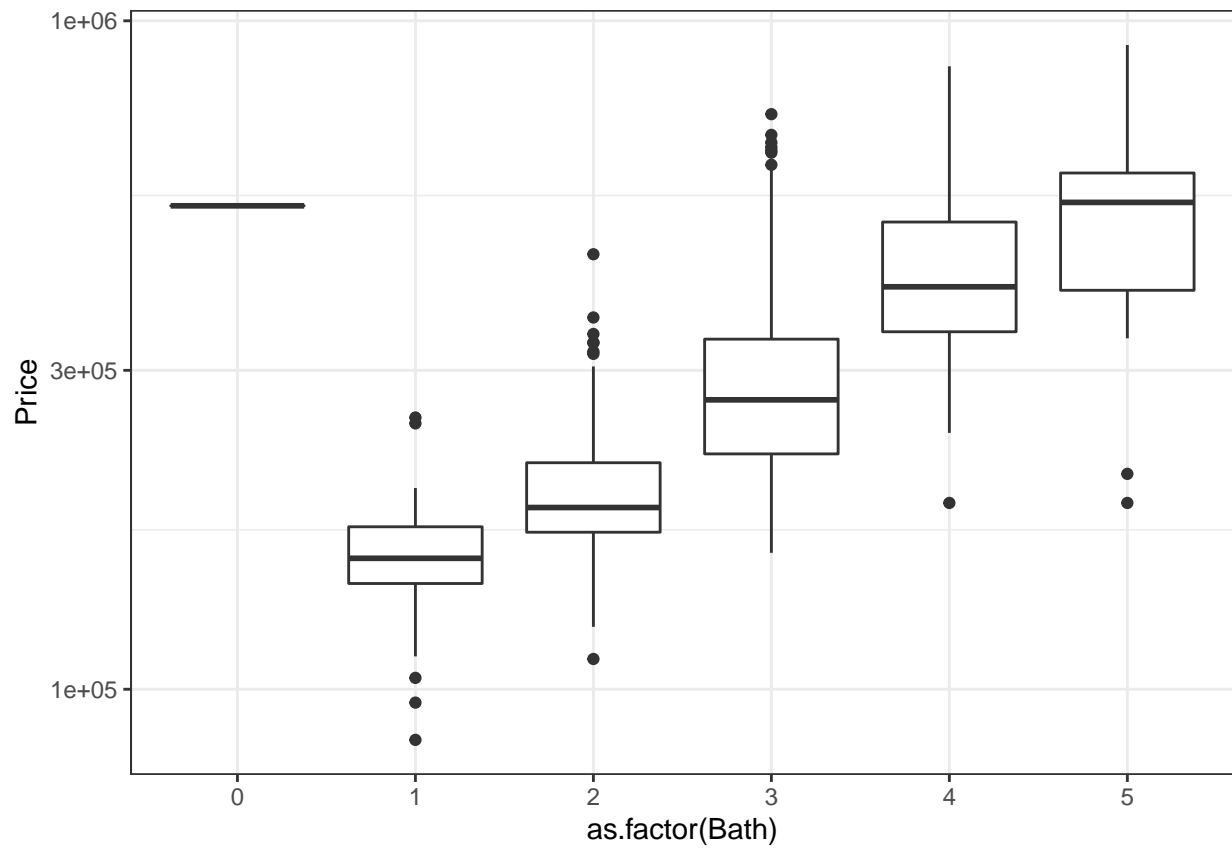


```
estate %>%
  mutate(Bath = recode(Bath, `6` = 5, `7` = 5)) ->
  estate
```

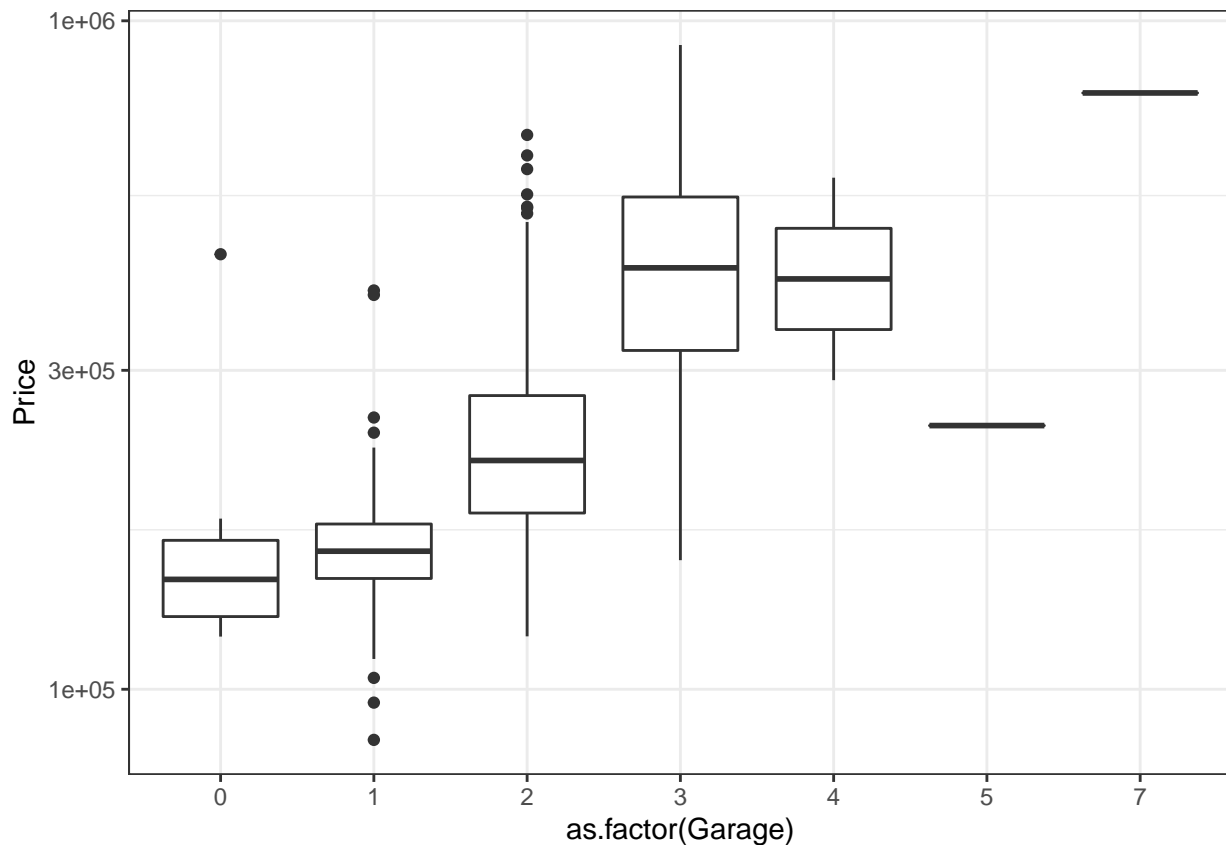
```
table(estate$Bath)
```

```
##
##  0  1  2  3  4  5
##  1 71 171 175 84 20
```

```
estate %>%
  ggplot(aes(x = as.factor(Bath), y = Price)) +
  geom_boxplot() +
  scale_y_log10()
```



```
## Garage also looks pretty linear
estate %>%
  ggplot(aes(x = as.factor(Garage), y = Price)) +
  geom_boxplot() +
  scale_y_log10()
```



Summary of interesting Observations

- There seems to be a bathroom saturation effect. Having more than 5 doesn't help you that much.
- Most quantitative relationships with price are exponential in nature. In other words, a multiplicative difference in area corresponds to a multiplicative difference in price. A multiplicative difference in age corresponds to a multiplicative difference in price.
- A lot of the discrete quantities (number of garages, number of baths, number of beds) can be treated as a linear relationship with log Price. So if you add one more bedroom, it results in some multiplicative change in price.
- There is one observation with no bed, no bath, but is super pricey. My guess is that this is an empty lot on some prime real estate. I would exclude this from any model and state clearly that our analysis is for homes with at least one bedroom.

Some exploratory forward/backward linear model stuff

```
estate %>%
  filter(Bath != 0) ->
  estate_sub
lmfull <- lm(logPrice ~ Bed + Bath + AC + Garage + Pool +
             Quality + Style + Highway + logArea + logLot + logAge,
             data = estate_sub)

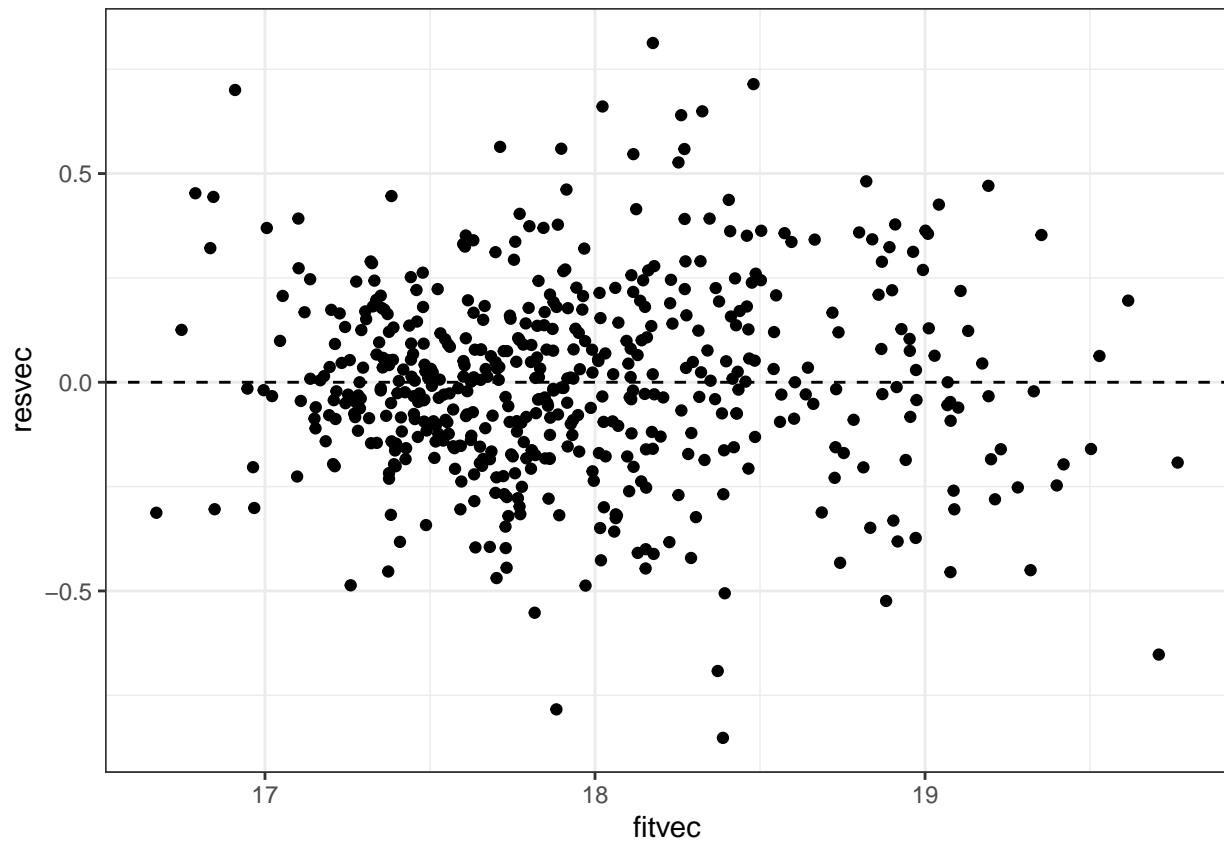
sout <- step(lmfull)

## Start:  AIC=-1487.72
## logPrice ~ Bed + Bath + AC + Garage + Pool + Quality + Style +
## Highway + logArea + logLot + logAge
```

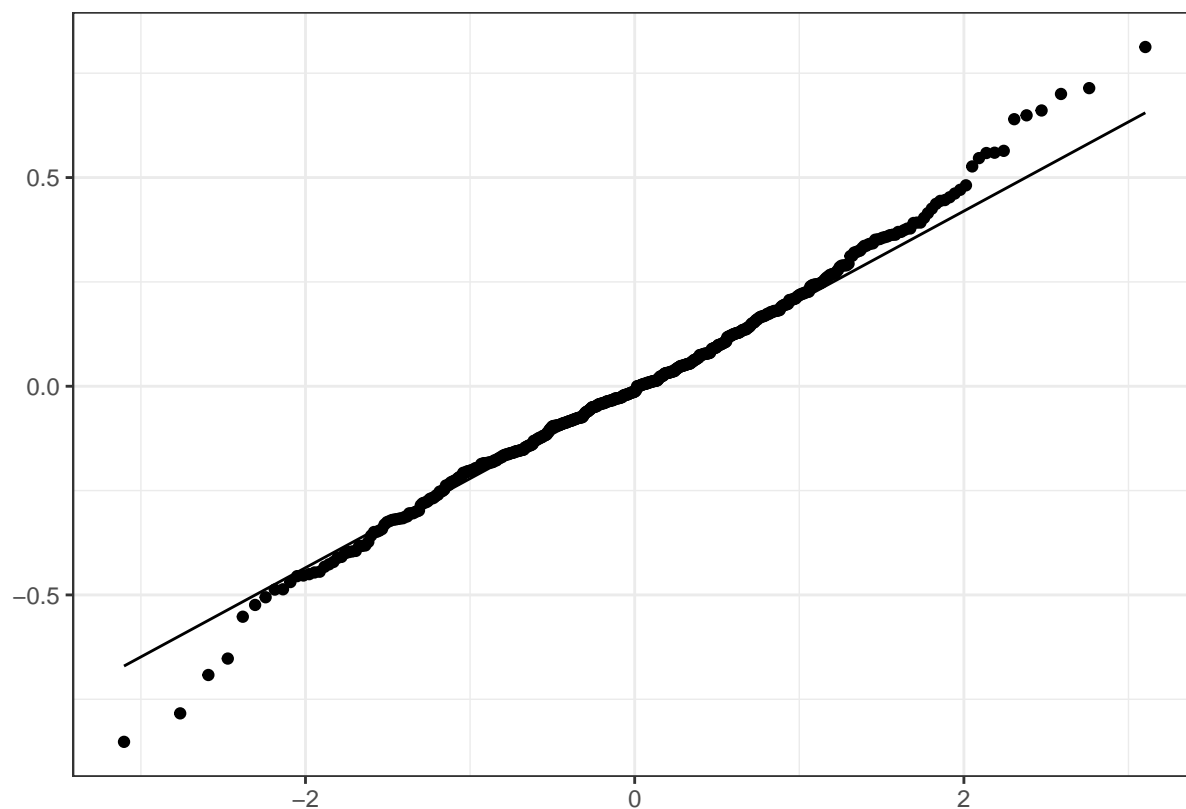
```
##
##           Df Sum of Sq   RSS   AIC
## - Bed      1    0.0663 27.716 -1488.5
## <none>                27.650 -1487.7
## - AC       1    0.1225 27.773 -1487.4
## - Garage   1    0.1489 27.799 -1486.9
## - Highway  1    0.2258 27.876 -1485.5
## - Pool     1    0.4227 28.073 -1481.8
## - Bath     1    0.8004 28.451 -1474.8
## - Style    9    2.0017 29.652 -1469.3
## - logLot   1    2.1256 29.776 -1451.1
## - Quality  2    3.7383 31.388 -1425.7
## - logAge   1    4.0607 31.711 -1418.3
## - logArea  1   10.4221 38.072 -1323.1
##
## Step:  AIC=-1488.47
## logPrice ~ Bath + AC + Garage + Pool + Quality + Style + Highway +
##           logArea + logLot + logAge
##
##           Df Sum of Sq   RSS   AIC
## <none>                27.716 -1488.5
## - AC       1    0.1389 27.855 -1487.9
## - Garage   1    0.1550 27.871 -1487.6
## - Highway  1    0.2196 27.936 -1486.4
## - Pool     1    0.4189 28.135 -1482.7
## - Bath     1    0.9985 28.715 -1472.0
## - Style    9    1.9839 29.700 -1470.5
## - logLot   1    2.1552 29.872 -1451.5
## - Quality  2    3.6775 31.394 -1427.6
## - logAge   1    3.9954 31.712 -1420.3
## - logArea  1   11.4289 39.145 -1310.6

resvec <- resid(sout)
fitvec <- fitted(sout)

## Residuals look pretty awesome
qplot(fitvec, resvec) +
  geom_hline(yintercept = 0, linetype = "dashed")
```



```
qplot(sample = resvec, geom = "qq") +  
  geom_qq_line()
```



```
coefvec <- coef(sout)
confintmat <- confint(sout)
sumlm <- summary(sout)
sumlm
```

```
##
## Call:
## lm(formula = logPrice ~ Bath + AC + Garage + Pool + Quality +
##     Style + Highway + logArea + logLot + logAge, data = estate_sub)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.85200 -0.15170 -0.01281  0.13651  0.81279
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    8.48969    0.60325   14.073 < 2e-16 ***
## Bath           0.07646    0.01800    4.248 2.57e-05 ***
## ACnoAC        -0.05112    0.03226   -1.584  0.11372
## Garage         0.03420    0.02043    1.674  0.09481 .
## PoolPool       0.11800    0.04288    2.752  0.00614 **
## QualityLow    -0.43037    0.05931   -7.256 1.53e-12 ***
## QualityMedium -0.34679    0.04286   -8.092 4.47e-15 ***
## Style2        -0.09211    0.03728   -2.471  0.01381 *
## Style3        -0.02283    0.03558   -0.641  0.52151
## Style4         0.09241    0.07399    1.249  0.21228
## Style5        -0.10048    0.06096   -1.648  0.09991 .
## Style6        -0.03690    0.06154   -0.600  0.54904
## Style7        -0.16793    0.03541   -4.743 2.75e-06 ***
```



```
## Style9          -0.15627    0.23832  -0.656  0.51230
## Style10         -0.33196    0.24201  -1.372  0.17078
## Style11         -0.53608    0.23684  -2.263  0.02403 *
## HighwaynoHighway 0.14513    0.07285   1.992  0.04689 *
## logArea         0.75983    0.05286  14.373 < 2e-16 ***
## logLot          0.11851    0.01899   6.242 9.23e-10 ***
## logAge          -0.14378    0.01692  -8.498 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2352 on 501 degrees of freedom
## Multiple R-squared:  0.862, Adjusted R-squared:  0.8568
## F-statistic: 164.8 on 19 and 501 DF, p-value: < 2.2e-16
```

Here is some example statements:

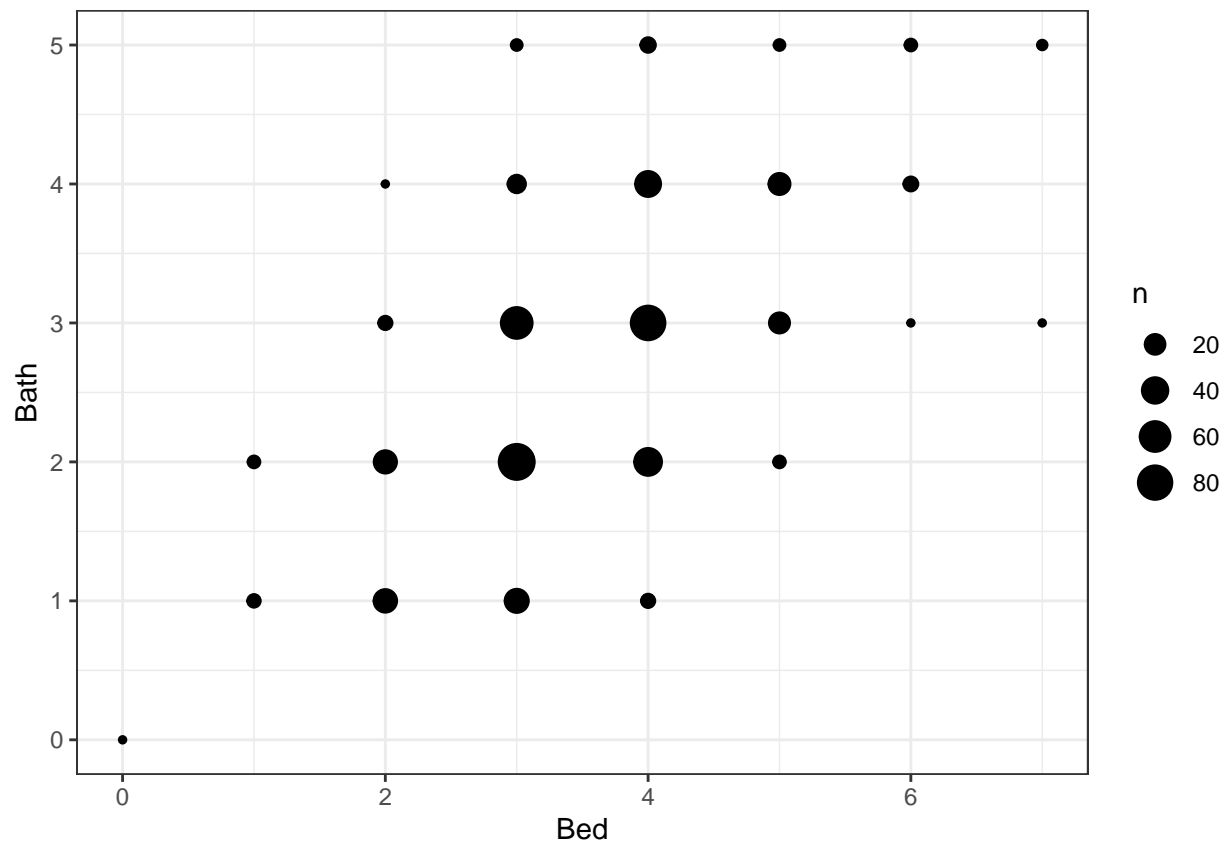
- Multiplicative increase in price when you add one bathroom: $r^{2^{\text{coefvec["Bath"]}}}$ (because I used base 2 when logging), which is 1.054 (95% CI $r^{2^{\text{confintmat["Bath",]}}$, which is 1.029 to 1.081).
- Doubling the area corresponds to a multiplicative increase of $r^{2^{\text{coefvec["logArea"]}}}$ (which is 1.693) in price (95% CI of $r^{2^{\text{confintmat["logArea",]}}$, which is 1.58 to 1.82).

I used causal language here, this is a little loose and relaxed. In formal write ups you should use non-causal language like:

- Houses with one more bathroom tend to cost 5% more (95% confidence of 3% to 8%)
- Homes with twice the area tend to cost 69% more (95% confidence of 58% to 82%).

Interestingly, bed was not informative given the other variables. Why? Probably because it is so highly correlated with the other variables that it doesn't add any additional information on price (at least given this dataset).

```
ggplot(estate, aes(x = Bed, y = Bath)) +
  geom_count()
```



```
ggplot(estate, aes(x = logArea, y = logPrice, color = as.factor(Bed))) +  
  geom_point()
```

