# Modeling House Prices Using Realtor Data

Agam Singh, Renos Zabounidis,
Manav Kulshrestha
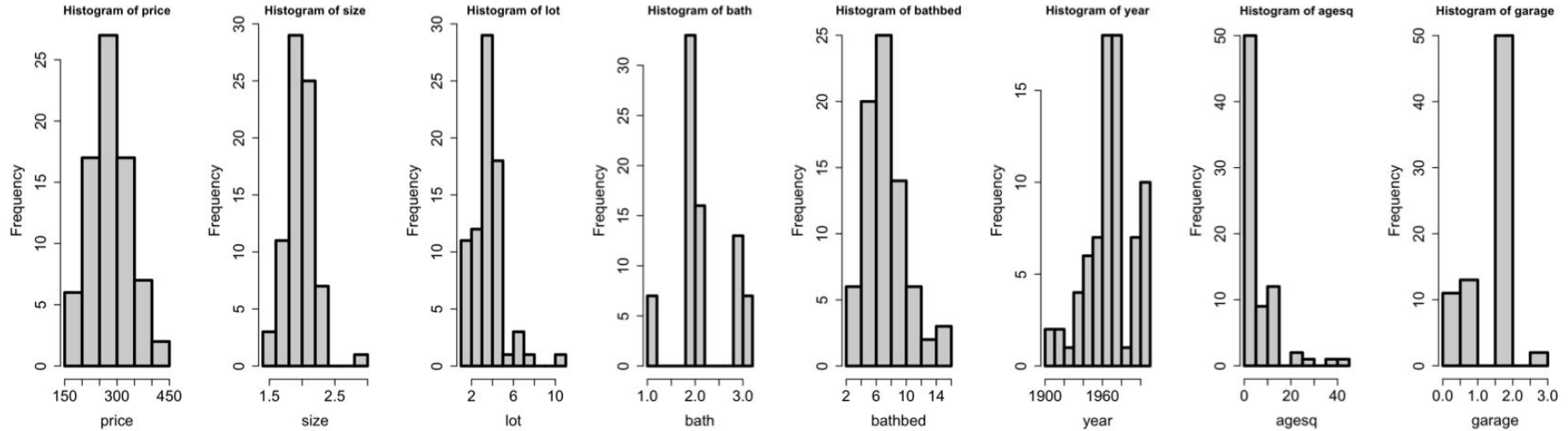and William Munson

# A Little Background Information

- Single-family homes in Eugene, Oregon
- Collected in 2005
- Mix of quantitative and qualitative variables
- 76 data points in total
- 18 variables in total (1 response and 17 predictors)
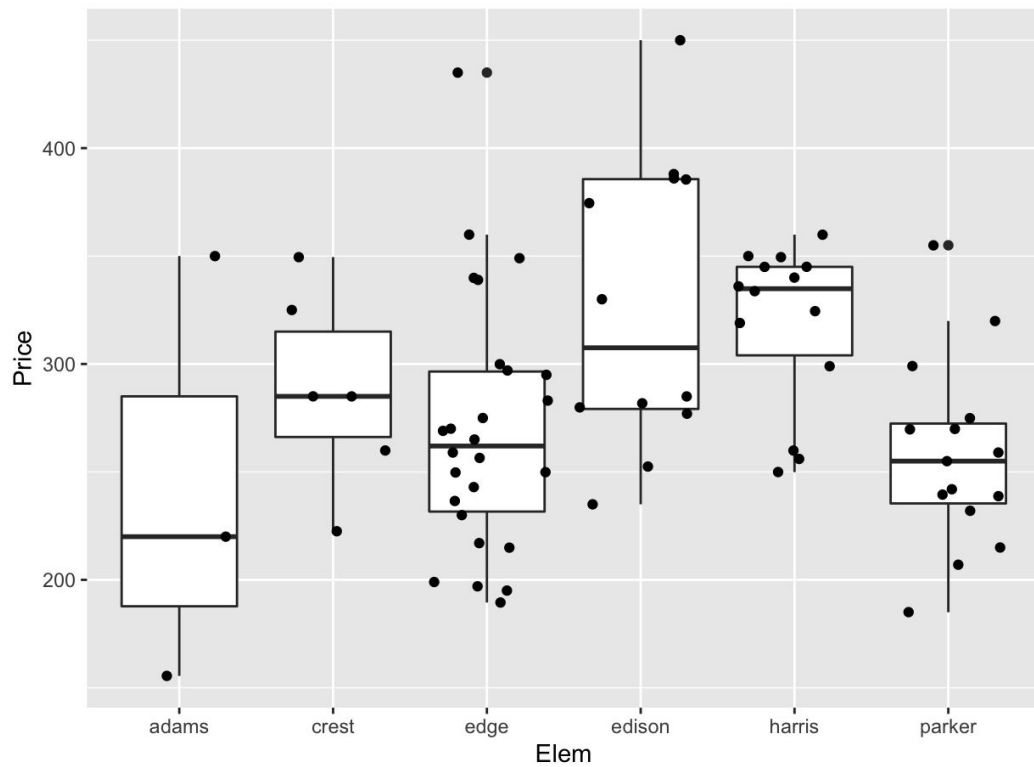
# A Closer Look at the Variables

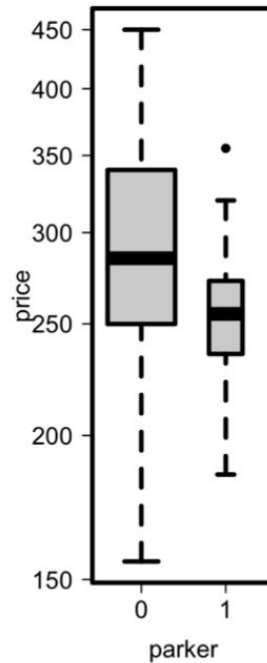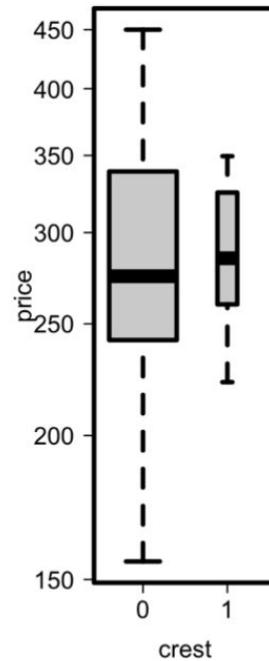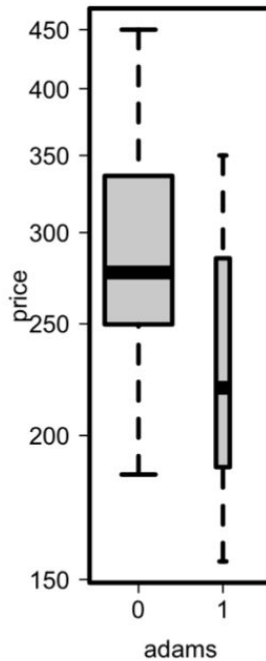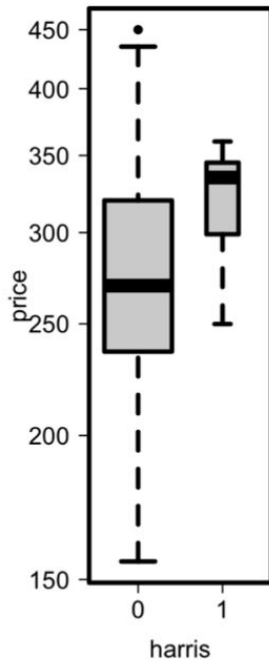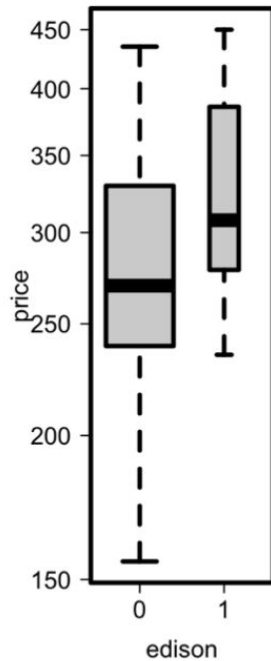| | |
|---:|:---|
| **Price** | `sale price (thousands of dollars)` |
| **Size** | floor size (thousands of square feet) |
| **Lot** | lot size category (from 1 to 11) |
| **Bath** | number of bathrooms (with half-bathrooms counting as 0.1) |
| **Bed** | number of bedrooms (between 2 and 6) |
| **BathBed** | interaction of Bath times Bed |
| **Year** | year built |
| **Age** | age (standardized: (Year-1970)/10) |
| **Agesq** | Age squared |
| **Garage** | garage size (0, 1, 2, or 3 cars) |
| **Status** | act (active listing), pen (pending sale), or sld (sold) |
| **Active** | indicator for active listing (reference: pending or sold) |
| **Elem** | nearest elementary school (edgewood, edison, harris, adams, crest, or parker) |
| **Edison** | indicator for Edison Elementary (reference: Edgewood Elementary) |
| **Harris** | indicator for Harris Elementary (reference: Edgewood Elementary) |
| **Adams** | indicator for Adams Elementary (reference: Edgewood Elementary) |
| **Crest** | indicator for Crest Elementary (reference: Edgewood Elementary) |
| **Parker** | indicator for Parker Elementary (reference: Edgewood Elementary) |

# The Quantitative Variables



They all seem approximately normal..

# The Schools

# The Schools

# Correlation Matrix

# Residual Plots

# Residual Plots

# Normal Plots

# Normal Plots

# Interaction Terms

- We included two additional
  interaction terms
  - bed*bath
  - edison*size

# Transformations

- We included a single transform, age^2

# Model Selection

# The Model

```
price = b0 + b1*size + b2*lot + b3*bath + b4*bed +
            b5*bathbed + b6*agesq + b7*garage +
     b8*active + b9*edison + b10*harris + b11*edisonsize


price = 355.151 + 43.773*size + 8.630*lot - 96.426*bath -
          78.393*bed + 29.347*bathbed + 1.409*agesq +
        20.203*garage + 29.040*active - 149.353*edison +
              50.222*harris + 110.178*edisonsize
```

# Model Fit

```
Residuals:
    Min      1Q  Median      3Q     Max
-106.97  -20.91    0.87   23.73   96.69

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept)   355.151    106.441   3.337 0.001416 **
size           43.773     28.981   1.510 0.135861
lot             8.630      3.364   2.565 0.012656 *
bath          -96.426     41.831  -2.305 0.024412 *
bed           -78.393     27.372  -2.864 0.005651 **
bathbed        29.347     11.688   2.511 0.014575 *
agesq           1.409      0.702   2.008 0.048909 *
garage         20.203      7.973   2.534 0.013741 *
active         29.040     10.860   2.674 0.009499 **
edison       -149.353    124.900  -1.196 0.236192
harris         50.222     14.418   3.483 0.000898 ***
edisonsize    110.178     64.593   1.706 0.092907 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 41.33 on 64 degrees of freedom
Multiple R-squared:  0.5995,	Adjusted R-squared:  0.5307
F-statistic:  8.71 on 11 and 64 DF,  p-value: 3.445e-09
```
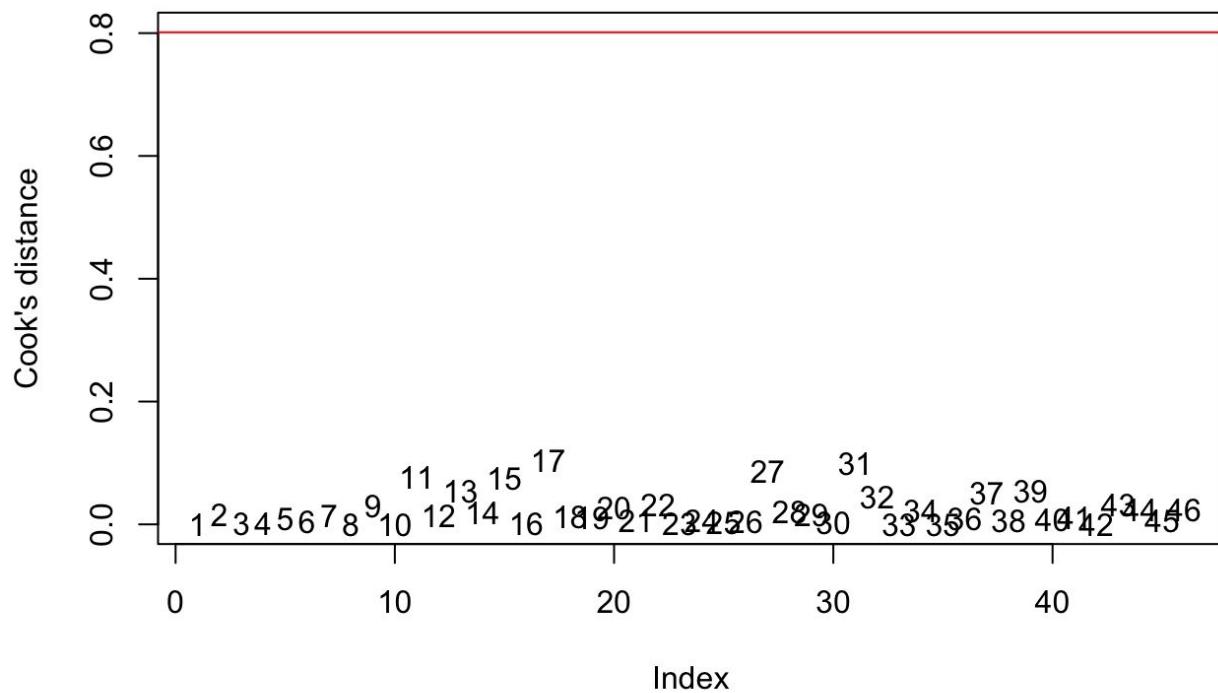
# Cross Validation

- 70/30 train/validation split

- MSPR: 1611.11

- MSE: 1438.56

- → *MSE* for the selected regression model is not seriously biased

- Achieved minimum MSE from all tested models

# No Outliers!

# Limitations

- We have too little data, there are only 76 data points
- There are only 53 points to fit to with a 70% training split
- Data only representative of Eugene, Oregon
- Data seems to have been pre-processed