1/28/2015

**Last time:**

- Encompassing idea of stat models partition variation
- $Y_i = \hat{y}_I + e_i$
- $Y_i$ = fitted value + residual
- = systematic + unpredictable
- Std deviation of $\hat{y}_I$ compared to $e_i$ to decide within or between group variation is the dominate model
- Basline/ offset form or dummy variable
  - Data in two groups ex) architecture and business
  - $X_i$ (group 1) is a dummy variable
    - $X_i = 1$ if case I is in group 1
    - $X_i = 0$ if case I is in group 0
  - $Y_i = B_o + B_1 X_i + e_i$
    - $B_1$ = difference between group means
  - If $X_i = 0$; $Y_i = B_o + e_i$
  - If $X_i = 1$; $Y_i = B_o + B_1 + e_i$
    - $X_i$ is the dummy variable
    - $B_o$ is the baseline/ intercept
    - $B_1$ is the slope or the coefficient on the dummy variable
  - More than 2 groups? More dummy variables
    - Groups: 0,1,2
    - $X_i = 1$ if case I is in group 1
    - $X_i = 2$ if case I is in group 2
    - $X_i = 0$ if case I is in group 0
    - Group 0 $Y_i = B_o + e_i$
    - Group 1 $y_i = B_o + B_1 + e_i$
    - Group 2 $y_i = B_o + B_2 + e_i$

**SAT R Script**

- UT200.csv load into Rstudio
- `mean(SAT.Q ~ School, data=ut2000)`
  - change mean to LM for get it on baseline offset form

  - `lm1 = lm(SAT.Q ~ School, data=ut2000)`

  - `coef(lm1)`

**OLS: ordinary least squares**

- $Y_i$ = response
  - Ex asking price of a truck on craigslist
  - $X_i$ is the predictor, feature, covariant, independent variable
  - Straight line: $Y_i = B_o + B_1 X_1$

- But this doesn't hit all the points
- Error is the "misses" of the equation
- Legendre says minimize the sum of $e_i^2$
  - $SUM(e_i^2,1,n)$
  - $SUM((Y_i-\hat{Y}_i)^2,1,n)$
  - $SUM(Y_i-(B_o+B_1x_1))^2$
- Why square?
  - No negatives
  - Square the errors punishes the outliers more (proportionate)
  - Real reason: Legendre did this by hand and to minimize the function you take the derivative of the square-> can't take the derivative of the absolute value
  - Deep conection between sum of sq errors and normal distributions *later in the course*
  - Variance deomposision: special property that can be invoked by pathagreom *we will look at this on Monday*

## Pickup walk through

- Upload dataset
  - Picture of coverage interval
  - (lm(price~miles, data=pickup)
    - Calculus in a micro second
- Save the output
  - Store the command as a model
    - Model1=lm(price~miles, data=pickup)
    - Find coef(model1)
    - Find fitted(model1)
    - Resid(model1)
    - Fitted(model1) + resid(model1)=original
  - Line of best fir
    - Plot(price~miles, data=pickup)
    - Abline(model1)
- 4 stories of linear regression
  1. Plug in prediction: plug in your value and the prediction comes out
     a. Sell truck on craigslist
     b. Newx = c(25000,50000,100000)
     c. Yhat= 14419.3762 + (-0.0643)*newx
  2. Summarize the trend: why we collect data
     a. What will the change in Y be if I wait to sell my truck online
     b. $\hat{Y}=B_o + B_1X$
        i. Take derivate with respect to X
        ii. $D\hat{y}/dx=B_1$
     c. Coef(model1) -6.42E-02 * 10000                    =-642.9944
  3. Taking the X-ness out of Y (statistical adjustment)

a. Residual can be thought of a Y adjusted for X
    4. Qualifying the Reduction in Uncertainty
        a. Sd(pickup$price)
        b. Guess sample mean of price, use std dev of prices to see how good that guess will be, graph and pick the point on the graph you are interested in
            i. Quantify with residuals
        c. Sd(resid(model1))
        d. Truecar.com makes money on data analysis
            i. Graph of price comparisions (stories 1-4)
                1. Normal distribution describing residuals on their stat model

Most important points: the uncertainty that remains after those 4 things are done is a huge factor to decision making

**Case Study:**

1) Define new variables that compute a mean value for each state's statistics (vmiles, mrall, and perinc) across all years. Remember how to compute groupwise means in R.

```
mean(mrall~statename, data=traffic2)
```
way more than 50 dots: we need to store means

mrall_mean = mean(mrall ~ statename, data traffic2)
vmiles_mean = mean(vmiles ~ statename, data=traffic2)

plot(mrall_mean~vmiles_mean, data=traffic2)
coef(fred)
abline(fred)

2) Make a scatter plot and fit a linear models that shows the relationship between fatality rate (response variable) and miles per driver (predictor).
3) Do the same thing, except with per capita income as the predictor.
4) Make a lattice plot that stratifies the scatter plot of fatality rate versus miles driven by the categorical variable indicating whether the state has mandatory jail for drunk driving. What does this plot suggest?