# Linear Regression

1. Overview and Basic Theory

2. R stuff and problems.

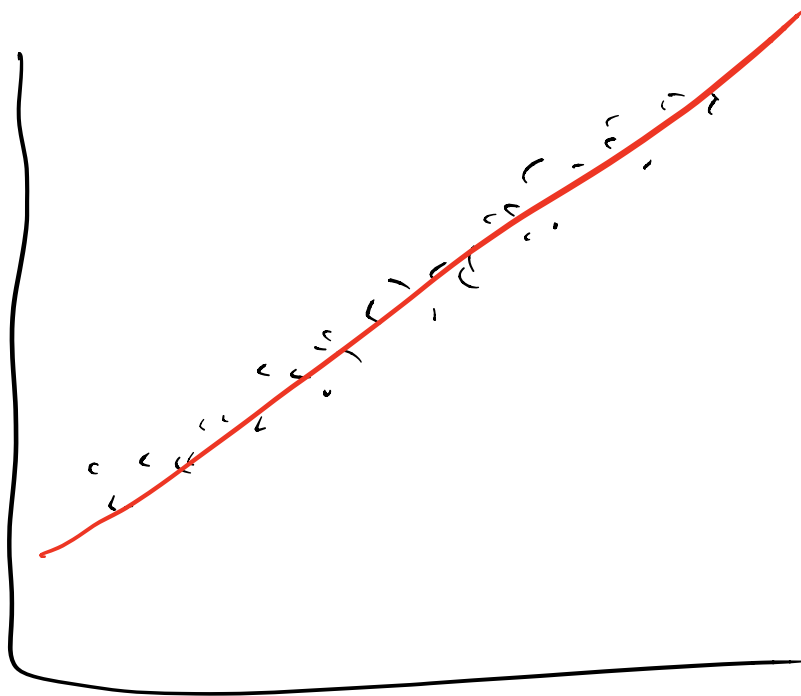   OLS

   SLR

   Least Squares

   Ld regression

   } Linear regression w/ Least Squares

# Linear Regression

We want to fit a straight line through data.



given $(x_i, y_i)$ data pairs.

Setup:

Assume we have a relationship

$$y_i \sim \beta_0 + \beta_1 x_i$$

response ← (under $y_i$)

covariate ← (under $x_i$)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad \leftarrow \text{noise}$$

$$\varepsilon_i \overset{iid}{\sim} [0, \sigma_\varepsilon^2]$$

mean    variance

How do we find

$\hat{\beta}_0$, $\hat{\beta}_1$ (estimates of the coefficients)

our line estimate becomes

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$ "fitted value"

we can minimize how bad we do.

$y_i - \hat{y}_i$ is how bad we do.

A simple choice
"Least squares"

$$\left(\hat{\beta_0}, \hat{\beta_1}\right) = \underset{b_0, b_1}{\arg\min} \sum_{i=1}^{n} \left(y_i - \hat{y}_i\right)^2$$

$$= \underset{b_0, b_1}{\arg\min} \sum_{i=1}^{n} \left(y_i - \left(b_0 + b_1 x_i\right)\right)^2$$

$$\Downarrow$$

$$\hat{\beta_1} = \frac{S_{xy}}{S_{xx}} = \frac{\sum_{i=1}^{n} \left(x_i - \bar{x}\right)\left(y_i - \bar{y}\right)}{\sum_{i=1}^{n} \left(x_i - \bar{x}\right)^2}$$

$$\hat{\beta_0} = \bar{y} - \hat{\beta_1}\bar{x}$$

again the setup

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$\varepsilon_i \overset{iid}{\sim} [0, \sigma_\varepsilon^2]$$

often we assume

$$\varepsilon_i \overset{ind}{\sim} N(0, \sigma_\varepsilon^2)$$

$$y_i \overset{ind}{\sim} N(\beta_0 + \beta_1 x_i, \sigma_\varepsilon^2)$$
$$i = 1, \dots, n$$

In our traditional

$$x_i \sim N(0, 1) \quad, i = 1, \dots, n$$

# The connection between Least squares and Normal MLE.

$$\left(\widehat{\beta_0}, \widehat{\beta_1}\right) = \underset{b_0, b_1}{\arg\min} \sum_{i=1}^{n} \left(y_i - (b_0 + b_1 x_i)\right)^2$$

MLE of $\beta_0, \beta_1$

$$y_i \overset{ind}{\sim} N\left(\beta_0 + \beta_1 x_i, \; \sigma_\varepsilon^2\right)$$

$$L(\beta_0, \beta_1) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma_\varepsilon^2}} \exp\left(-\frac{1}{2}\left(\frac{\left(y_i - (\beta_0 + \beta_1 x_i)\right)^2}{\sigma_\varepsilon^2}\right)\right)$$

$$= \left( \frac{1}{\sqrt{2\pi\sigma_\varepsilon^2}} \right)^n \exp\left( \frac{-1}{2\sigma_\varepsilon^2} \sum_{i=1}^{n} \left( y_i - (\beta_0 + \beta_1 x_i) \right)^2 \right)$$

$$\ell_n(\beta_0, \beta_1) = \log\left( \left( \frac{1}{\sqrt{2\pi\sigma_\varepsilon^2}} \right)^n \right)$$

$$- \frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^{n} \left( y_i - (\beta_0 + \beta_1 x_i) \right)^2$$

equiv to maximizing

$$- \frac{1}{2\sigma_\varepsilon^2} \sum_{i=1}^{n} \left( y_i - (\beta_0 + \beta_1 x_i) \right)^2$$

equiv to minimizing

$$\sum_{i=1}^{n} \left( y_i - \left( \beta_0 + \beta_1 x_i \right) \right)^2$$

## Correlation

For random variables

$$\rho = cor(X, Y)$$

$$= \frac{Cov(X, Y)}{sd(x) \cdot sd(y)}$$

$$r = \frac{\sum\limits_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum\limits_{i=1}^{n} (x_i - \bar{x})^2} \; \sqrt{\sum\limits_{i=1}^{n} (y_i - \bar{y})^2}}$$

$$\hat{\beta}_1 = r_{xy} \frac{sd_y}{sd_x}$$

$$\underset{b_0, b_1}{\operatorname{argmin}} \sum_{i=1}^{n} \left| y_i - (b_0 + b_1 x_i) \right| \quad \text{(L1)}$$

$$\text{vs.}$$

$$\underset{b_0, b_1}{\operatorname{argmin}} \sum_{i=1}^{n} \left( y_i - (b_0 + b_1 x_i) \right)^2 \quad \text{(L2)}$$

## Residuals

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

$$\varepsilon_i \overset{iid}{\sim} N(0, \sigma_\varepsilon^2)$$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \text{"fitted values"}$$

Residuals

$$\hat{\varepsilon}_i = y_i - \hat{y}_i$$

$$\left( \text{Least squares} \atop \min \sum_{i=1}^{n} \hat{\varepsilon}_i^2 \right)$$