## CH 8: Least Squares and Maximum Likelihood

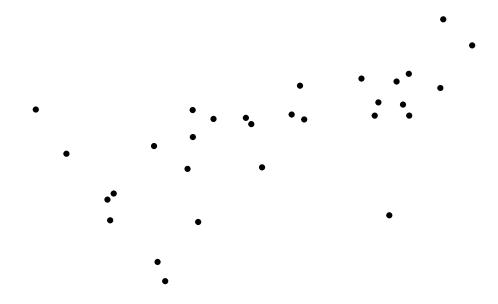
## Least Squares

Least squares is a common method for estimating regression coefficients in a linear model. Consider the model,  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ , where i indexes the  $i^{th}$  observation.

Then define a residual as:

$$r_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)$$

Least squares refers to minimizing the squared residuals (or the difference between the regression line and the points).



The residual sum of square are formally defined as:

$$RSS = \sum_{i=1}^{n} \left( y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2$$

Hence, to minimize this value, we will take the derivatives with respect to  $\beta_0$  and  $\beta_1$ 

$$\frac{dRSS}{d\beta_0} = -2\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)$$

$$(\text{set} = 0) - 2\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \beta_0 - \sum_{i=1}^n \beta_1 x_i = 0$$

$$n\beta_0 = \sum_{i=1}^n y_i - \sum_{i=1}^n \beta_1 x_i$$

$$\beta_0 = \frac{\sum_{i=1}^n y_i}{n} - \beta_1 \frac{\sum_{i=1}^n x_i}{n}$$

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

Then we also solve for  $\beta_1$  plugging in the value of  $\beta_0$ 

$$\frac{dRSS}{d\beta_{1}} = -2\sum_{i=1}^{n} (y_{i} - \beta_{0} - \beta_{1}x_{i}) x_{i}$$

$$(\text{set} = 0) - 2\sum_{i=1}^{n} (y_{i} - \beta_{0} - \beta_{1}x_{i}) x_{i} = 0$$

$$\sum_{i=1}^{n} (y_{i} - \beta_{0} - \beta_{1}x_{i}) x_{i} = 0$$

$$\sum_{i=1}^{n} x_{i}y_{i} - \sum_{i=1}^{n} x_{i}(\bar{y} - \beta_{1}\bar{x}) - \sum_{i=1}^{n} \beta_{1}x_{i}^{2} = 0$$

$$\sum_{i=1}^{n} x_{i}y_{i} - \sum_{i=1}^{n} x_{i}(\bar{y} - \beta_{1}\bar{x}) - \sum_{i=1}^{n} \beta_{1}x_{i}^{2} = 0$$

$$\beta_{1} \left(\sum_{i=1}^{n} x_{i}^{2} - \bar{x}\sum_{i=1}^{n} x_{i}\right) = \sum_{i=1}^{n} x_{i}y_{i} - \bar{y}\sum_{i=1}^{n} x_{i}$$

$$\beta_{1} \left(\sum_{i=1}^{n} x_{i}^{2} - \bar{x}\sum_{i=1}^{n} x_{i}\right) = \sum_{i=1}^{n} x_{i}y_{i} - \sum_{i=1}^{n} \frac{y_{i}}{n}n\bar{x}$$

$$\beta_{1} \left(\sum_{i=1}^{n} x_{i}^{2} - \bar{x}\sum_{i=1}^{n} x_{i}\right) = \sum_{i=1}^{n} (x_{i} - \bar{x}) y_{i}$$

$$\beta_{1} \left(\sum_{i=1}^{n} x_{i}^{2} - n\bar{x}^{2}\right) = \sum_{i=1}^{n} (x_{i} - \bar{x}) y_{i}$$

$$\vdots = \sum_{i=1}^{n} (x_{i} - \bar{x}) y_{i}$$

$$\vdots = \sum_{i=1}^{n} (x_{i} - \bar{x}) y_{i}$$

$$\beta_{1} \sum_{i=1}^{n} (x_{i} - \bar{x}) y_{i}$$

This calculation is actually considerably easier using matrix notation.

Let the model be specified as  $\mathbf{y} = X\beta + \epsilon$ , where

•  $\mathbf{y}$  is a  $n \times 1$  matrix (or vector), such that

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

• X is a  $n \times 2$  matrix,

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

- $\beta$  is  $2 \times 1$  matrix,  $\beta = (\beta_0 \ \beta_1)$
- $\epsilon$  is a  $n \times 1$  matrix (or vector), such that

$$\epsilon = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

To verify this matrix algebra, we can look at the first row of the matrix which results in

$$y_1 = \beta_0 + \beta_1 x_1 + \epsilon_1$$

The residuals can be defined as

$$\mathbf{r} = \mathbf{y} - X\hat{\beta} = \begin{pmatrix} y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1 \\ y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2 \\ \vdots \\ y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n \end{pmatrix}$$

Then the sum of squares is written as:

$$\sum_{i=1}^{n} = \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i\right)^2 = \mathbf{r}^T \mathbf{r}$$

To minimize the sum of squares of the residuals, consider

$$\mathbf{r}^{T}\mathbf{r} = (\mathbf{y} - X\beta)^{T}(\mathbf{y} - X\beta)$$

$$= \mathbf{y}^{T}\mathbf{y} - \mathbf{y}TX\beta - \beta^{T}X^{T}\mathbf{y} + \beta^{T}X^{T}X\beta$$

$$= \mathbf{y}^{T}\mathbf{y} - 2\beta^{T}X^{T}\mathbf{y} + \beta^{T}X^{T}X\beta$$

Now take the derivative

$$\frac{d \, SSR}{d \, \beta^T} = -2X^t \mathbf{y} + 2X^T X \beta \tag{1}$$

and set = 0

$$0 - 2X^T \mathbf{y} + 2X^T X \beta \tag{2}$$

$$X^T \mathbf{y} = X^T X \beta \tag{3}$$

$$(X^T X)^{-1} X^T \mathbf{y} = \beta \tag{4}$$

 $\sigma$  is also estimated with the sum of the squared residuals

$$\hat{\sigma} = \sqrt{\frac{1}{n-k} \sum_{i=1}^{n} (y_i - X_i \hat{\beta})^2}$$

where k is the number of predictors in the model (including the intercept)

With least-squares estimation, there is no specification of a probability distribution. This is strictly a geometric procedure. Nevertheless, the result is the same as what is known as the maximum likelihood estimates.

Recall a regression model can be written as

$$y = \beta_0 + \beta_1 x + \epsilon, \ \epsilon \sim N(0, \sigma^2)$$

or

$$y \sim N(\beta_0 + \beta_1 x, \sigma^2)$$

Hence, the likelihood corresponds to how well the data point y corresponds to the normal density (or likelihood),

$$N(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2\right)$$

## lm and stan\_glm

With Bayesian inference, and the use of prior information, the posterior is a product of the likelihood of the data as well as the prior.

The role prior distribution is often characterized as a penalty (to the likelihood) or regularization, where the prior can down weight some values of the parameter.

As we have seen, Bayesian inference enables simulated based approaches for summarizing parameter coefficients, contrasts, and predictions. Bayesian inference also allows uncertainty to be expressed using probability, as opposed to using confidence. However, the cost is specifying a prior, which can be subjective.

Given that maximum-likelihood methods are optimizing the likelihood while Bayesian inference focus on the posterior (likelihood + prior), we'd expect differences in the results.

The default priors are weakly informative, so that posterior is not vastly different from the likelihood.

```
beer <- read_csv('http://math.montana.edu/ahoegh/Data/Brazil_cerveja.csv')</pre>
## Parsed with column specification:
## cols(
##
     consumed = col_double(),
     precip = col_double(),
##
    max_tmp = col_double(),
##
     weekend = col double()
## )
beer %>% lm(consumed ~ max_tmp, data = .) %>% coef()
## (Intercept)
                   max_tmp
     7.9749394
                 0.6548456
beer %>%
  stan_glm(consumed ~ max_tmp, data = ., refresh = 0, iter = 100000) %>% coef()
## (Intercept)
                   max_tmp
      7.982674
##
                  0.654593
```

We can explicitly state that flat (uniform, uniformative) priors, so that the posterior and the likelihood are are the same

```
beer <- read_csv('http://math.montana.edu/ahoegh/Data/Brazil_cerveja.csv')</pre>
## Parsed with column specification:
## cols(
##
     consumed = col_double(),
     precip = col_double(),
##
##
    max_tmp = col_double(),
##
     weekend = col_double()
## )
beer %>% lm(consumed ~ max_tmp, data = .) %>% coef()
## (Intercept)
                   max_tmp
    7.9749394
##
                 0.6548456
beer %>%
  stan_glm(consumed ~ max_tmp, data = .,
           refresh = 0, iter = 100000,
           prior_intercept = NULL,
           prior = NULL,
           prior_aux = NULL) %>% coef()
## (Intercept)
                   max_tmp
    7.9775519
                 0.6546636
```

Uncertainty Intervals
The authors suggest calli



I will be more lenient about this, but outside of my class, make sure you are clearly articulating the type of estimation and appropriate interpretation.