

## Multiple Linear Regression

Suppose our model is

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$$

and we observe  $y_1, \dots, y_n$ .

We can summarize all the equations as:

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

$$Y = X\beta + \varepsilon$$

The least squares approach would be to minimize  $\sum \varepsilon_i^2 = \varepsilon^T \varepsilon$ .

$$\text{Since } \varepsilon = Y - X\beta, \text{ we look to } (Y - X\beta)^T (Y - X\beta) = \\ = Y^T Y - \beta^T X^T Y - Y^T X \beta + \beta^T X^T X \beta$$

$$\text{Now } \frac{\partial}{\partial \beta} \varepsilon^T \varepsilon = \frac{\partial}{\partial \beta} (Y^T Y - \beta^T X^T Y - Y^T X \beta + \beta^T X^T X \beta) = \\ = -1^T X^T Y - Y^T X \cdot 1 - 1^T X^T X \beta + \beta^T X^T X \cdot 1 = -2 X^T Y + 2 X^T X \beta$$

$$\Rightarrow \text{when } X^T X \beta = X^T Y \text{ ie when } \hat{\beta} = (X^T X)^{-1} X^T Y.$$

In the simple linear regression case,

$$X^T X = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \quad (X^T X)^{-1} = \begin{bmatrix} \frac{\sum x_i^2}{n S_{xx}} & -\bar{x}/S_{xx} \\ -\bar{x}/S_{xx} & 1/S_{xx} \end{bmatrix} = \begin{bmatrix} c_{00} & c_{01} \\ c_{10} & c_{11} \end{bmatrix}$$

The estimator  $\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y$

From this we can calculate the

$$\begin{aligned} SSE &= \sum (y_i - \hat{y}_i)^2 = (Y - X\hat{\beta})^T (Y - X\hat{\beta}) = \\ &= Y^T Y - Y^T X\hat{\beta} - \hat{\beta}^T X^T Y + \hat{\beta}^T X^T X \hat{\beta} = \\ &= Y^T Y - 2\hat{\beta}^T X^T Y + \hat{\beta}^T X^T X (X^T X)^{-1} X^T Y = Y^T Y - 2\hat{\beta}^T X^T Y + \hat{\beta}^T \mathbf{1} \cdot X^T Y = Y^T Y - \hat{\beta}^T X^T Y. \end{aligned}$$

All our previous properties carry over for  $\hat{\beta}$ :

Assume  $\varepsilon_i$  independent.  $E=0$   $V=\sigma^2$

- 1  $E\hat{\beta}_i = \beta_i$
- 2  $V\hat{\beta}_i = c_{ii} \sigma^2$  where  $c_{ii} = (X^T X)^{-1}_{ii}$
- 3  $cov(\hat{\beta}_i, \hat{\beta}_j) = c_{ij} \sigma^2$  where  $c_{ij} = (X^T X)^{-1}_{ij}$
- 4  $SSE/(n-k-1)$  is unbiased for  $\sigma^2$ .
- 5 If  $\varepsilon_i$  are normally distributed
- 5 Each  $\hat{\beta}_i$  is normal
- 6  $SSE/\sigma^2 \sim \chi^2(n-k-1)$
- 7  $SSE$  and  $\hat{\beta}_i$  are independent

## Inferences on $a^T \beta$

Suppose that we want to make inferences on a linear combination of model parameters.

The linear combination  $a^T \beta$  has (by linearity of expectation) an unbiased estimator  $a^T \hat{\beta}$ . It is normally distributed if  $\epsilon \sim N(0, \sigma^2 I)$ .

By expanding  $V(a^T \hat{\beta})$ , we can recognize the result as

$$V(a^T \hat{\beta}) = [a^T (X^T X)^{-1} a] \sigma^2$$

With a normally distributed estimator with known variance we can construct Wald and T tests and CIs.

We get T-test test statistics

$$T = \frac{a^T \hat{\beta} - (a^T \beta)_0}{S \sqrt{a^T (X^T X)^{-1} a}} \sim t(n-k-1) \quad \text{where } S^2 = \frac{Y^T Y - \hat{\beta}^T X^T Y}{n-k-1}$$

$$a^T \beta \in a^T \hat{\beta} \pm t_{\alpha/2} S \sqrt{a^T (X^T X)^{-1} a} \quad \text{where } t_{\alpha/2} = F_{t(n-k-1)}^{-1}(\alpha)$$

This setup can be used both for inferences on each  $\beta_i$  in isolation, and also for estimation and prediction at a new value  $x^* = (x_1^*, \dots, x_k^*)$ .

For estimation, we'd just use the T-statistic as  $\hat{\epsilon}_i$ .

For prediction, we instead (again) estimate the error  $y^* - \hat{y}^*$

The same derivation as before leads us to the variance  $V_{\text{error}} = [1 + a^T(X^T X)^{-1} a] \sigma^2$  and corresponding CIs.

## The partial F-test

Suppose we want to figure out whether or not to include all variables when building a model.

This would be of interest eg to avoid overfitting where the model learns more about the data than about generalizing from data

We can do this by fitting linear models both with and without the extra variables and see if sufficiently much variance is explained to motivate the extra data used.

We would have a complete model

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_g x_g + \beta_{g+1} x_{g+1} + \dots + \beta_k x_k + \varepsilon$$

and a reduced model

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_g x_g + \varepsilon$$

With data we can calculate SSE for both, producing  $SSE_R$  and  $SSE_C$ .

We expect  $SSE_R \geq SSE_C$ , and so we can use their difference as a measure of the importance of  $x_{g+1}, \dots, x_k$ .

A hypothesis test would use  $H_0: \beta_{g+1} = \dots = \beta_k = 0$  vs  $H_A: \text{at least one } \neq 0$ .

We have, in the end, three sums of squares — each related to a  $\chi^2$  variable — much like in our ANOVA earlier. assuming  $H_0$ , the following holds:

$$\begin{aligned} \text{SSE}_R & \quad X_3 = \text{SSE}_R / \sigma^2 \sim \chi^2(n-g-1) \\ \text{SSE}_C & \quad X_2 = \text{SSE}_C / \sigma^2 \sim \chi^2(n-k-1) \\ \text{SSE}_R - \text{SSE}_C & \quad X_1 = (\text{SSE}_R - \text{SSE}_C) / \sigma^2 \sim \chi^2(k-g) \end{aligned}$$

$X_1$  and  $X_2$  are independent random variables.

If the extra variables are important, then the reduction of variance will be large — so  $X_1$  will be large.

With two  $\chi^2$  variables it is quite tempting to compute an F-statistic:

$$F = \frac{X_1 / (k-g)}{X_2 / (n-k-1)} \sim F_{\frac{k-g}{n-k-1}}$$

greater values of this F indicate higher importance for the extra variables

By default, R will compute this F-statistic for you with  $g=1$  — ie comparing

$$Y = \beta_0 + \varepsilon \quad \text{to} \quad Y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \varepsilon$$

You get details on the F calculation using `anova(model)`.

For other  $g$ , you will need to fit both models yourself and feed them to `anova`.

### Comparison with $Y = \beta_0 + \varepsilon$

In this model,  $SSE_R = S_{yy}$ .

We can compare the full variation (in  $S_{yy}$ ) to the variation after accounting for  $x_1, \dots, x_k$  (in  $SSE_c$ ) by comparing their difference  $S_{yy} - SSE_c$  to  $S_{yy}$ .

We define the multiple coefficient of determination

$$R^2 = \frac{S_{yy} - SSE_c}{S_{yy}}$$

$R^2$  measures proportion of variance explained by all the variables.

## ANOVA as a linear model

Recall that ANOVA tests  $H_0: \mu_1 = \dots = \mu_k = 0$  for a set of  $k$  groups.

This  $H_0$  is quite similar to the  $H_0: \beta_1 = \dots = \beta_k = 0$ .

No coincidence: Let  $x_{ij} = 1$  if observation  $i$  belongs to group  $j$ . Then  $Y = X\beta + \varepsilon$  reduces straight to the ANOVA model.

## Linear independence

If the columns in  $X$  are not linearly independent, then  $(X^T X)^{-1}$  does not exist and the least squares estimation fails.

In an ANOVA, this is guaranteed to be the case because of the presence of  $\beta_0$ .

The usual solution is to just drop a column (eg the last one) and use  $\beta_0$  instead.