# §27. Linear Regression

The model is

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

where

Y - dependent variable (or response)

 $x_i$  - independent variables (or predictors)

 $\beta_i$  - regression coefficients ,  $\beta_0$  - intercept

 $\varepsilon$  - noise term

# Least Squares Estimation of The Parameters

Say

$$(x_{i1}, x_{i2}, \dots, x_{ik}, y_i)$$
;  $i = 1, \dots, n, n > k$ 

are observed data points.

Let

$$\varepsilon_i = y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_n x_{ik})$$

be the residuals.

Then

$$L = \sum_{i=1}^{k} \varepsilon_i^2$$

is the sum of squared 'errors'. This is the quantity we would like to minimize:

$$\frac{\partial L}{\partial \beta_j} = -2\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^k \beta_j x_{1j} \right) x_{ij} = 0 \tag{1}$$

and we can solve for  $\beta_i$ . It is best to write everything in matrix notation

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \qquad X = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1k} \\ \vdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{bmatrix} \qquad \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} \qquad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Remark This model covers nonlinear cases as well. E.g.  $x_1 = x, \quad x_2 = x^2,$   $x_3 = xz, \quad x_4 = z,$   $x_5 = z^2$ So:  $Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 xz$  $+ \beta_4 z + \beta_5 z^2 + \varepsilon$  Then the solutions of the least squares equations (1) are

$$\hat{\beta} = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}y$$

The hat on  $\hat{\beta}$  signifies values estimated from the data. For the fitted model predictions, we have  $\hat{y} = X\hat{\beta}$  and the residuals are

 $e = y - \hat{y}$ 

Assume that the errors,  $\varepsilon_i$ , are independent and identically distributed distributed random variables with mean 0, and variance  $\sigma^2$ .

Then

1.  $\hat{\beta} = (X^{\mathsf{T}}X)^{-1}X^{\mathsf{T}}y$  is an unbiased estimator for  $\beta$ :

$$E(\hat{\beta}) = \beta$$

2. An unbiased estimator for  $\sigma^2$  is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n e_i^2}{n-p} = \frac{SS_E}{n-p}$$

where p = k + 1 is the number of parameters in the model.

### Significance of The Regression

Many different test for the 'quality' of the regression could be done. Here are some.

1. Tests on the individual repression coefficients:

 $H_0 = \beta_i = \beta_{j0}$  versus  $H_1 : \beta_j \neq \beta_{j0}$ 

The test statistic is:

$$t_0 = \frac{\hat{\beta}_j - \beta_{j0}}{\sqrt{\sigma^2 c_{jj}}} \qquad ; \qquad (n-p) \, \mathrm{df}$$

where  $c_{ij}$  is the  $j^{th}$  diagonal element of  $(X^{\mathsf{T}}X)^{-1}$ .

The most important special case is  $H_0: \beta_i = 0$  versus  $H_1 = \beta_i \neq 0$ and  $H_0$  is not rejected. This indicates that the regressor,  $x_i$ , can be deleted from the model.

2.  $R^2$  and  $R^2_{adi}$ 

$$SS_E = \sum_{i=1}^{m} \varepsilon_i^2;$$
  $SS_T = \sum_{i=1}^{n} y_i^2 - \frac{1}{n} \left(\sum y_i\right)^2;$   $SS_R = SS_T - SS_E$ 

$$R^2 = \frac{SS_R}{SS_T} = 1 - \frac{SS_E}{SS_T}$$

#### **Example**

 $R^2$  = 0.98 indicates that the model accounts for 98% of the variability in the data.

Note:  $R^2$  can never decrease when a regressor is added. Adding a new regressor which only marginally increases  $R^2$  could be counterproductive - the model is less interpretable; Occam's razor. Thus, the adjusted  $R^2$  statistic is (heavily) used:

$$R_{adj}^2 = 1 - \frac{SS_E/(n-p)}{SS_T/(n-1)}$$

 $R_{adj}^2$  penalizes the analyst for adding terms to the model and is an easy way to prevent over fitting, including regressors which are not really useful.  $R_{adj}^2$  is used for variable selection.

# Confidence and Prediction Intervals in Regression

A  $100(1 - \alpha)\%$  confidence interval on the **regression coefficient**  $\beta$ ; is

$$\hat{\beta}_j - t_{\alpha/2} \sqrt{\hat{\sigma}^2 c_{jj}} \leq \beta_i \leq \hat{\beta}_j + t_{\alpha/2} \sqrt{\hat{\sigma}^2 c_{jj}} \quad ; \quad (n-p) df$$

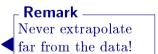
A  $100(1-\alpha)\%$  confidence interval on the **mean response** at the point  $x_{01}, \ldots, x_{0k}$  is

$$\hat{\mu}_{Y|x_0} - t_{\alpha/2} \sqrt{\hat{\sigma}^2 x_0^{\mathsf{T}} (x^{\mathsf{T}} x)^{-1} x_0} \leq \mu_{Y|x_0} \leq \hat{\mu}_{Y|x_0} + t_{\alpha/2} \sqrt{\hat{\sigma}^2 x_0^{\mathsf{T}} (x^{\mathsf{T}} x)^{-1} x_0}$$

A  $100(1-\alpha)\%$  prediction interval for a future observation of the response, Y, at  $x_{01}, \ldots, x_{0k}$  is

$$\hat{y}_0 - t_{\alpha/2} \sqrt{\sigma^2 (1 + x_0^{\mathsf{T}}(x^{\mathsf{T}}x)^{-1}x_0)} \leq Y_0 \leq \hat{y}_0 + t_{\alpha/2} \sqrt{\sigma^2 (1 + x_0^{\mathsf{T}}(x^{\mathsf{T}}x)^{-1}x_0)}$$

This prediction interval expresses both the error in estimating the mean of  $Y \mid x_0$  as well as the inherent variability of Y at fixed  $x = x_0$ .



# Selection of Variables and Model Building

- Step-wise regression. Adding or removing variables at each step based on F-test
- Forward selection: variables are added me at a time
- Backwards elimination: start will all possible regressors and eliminate the insignificant ones one at a time.

## **Multicollinearity**

It is expected that there might be dependencies between the regressors themselves.

Let  $R_j^2$  be the coefficient of determination resulting from regressing  $x_j$  on the remaining k-1 regressors. The variance of  $\hat{\beta}_j$  is effectively 'inflated' with variance inflation factor for  $\beta_j$ :

$$\text{VIF}\left(\beta_{j}\right) = \frac{1}{1 - R_{j}^{2}}$$

Estimation of the regression coefficients is very imprecise when multicollinearity is present. To combat multicollinearity it might be possible to collect more data or simply use a different (nonlinear regression) model.