

## Lecture 12. The linear model

# Introduction to linear models

- Linear models can be used to explain or model the relationship between a *response*, or *dependent*, variable and one or more *explanatory variables*, or *covariates* or *predictors*.
- For example, how do motor insurance claims depend on the age and sex of the driver, and where they live?  
Here the claim rate is the response, and age, sex and region are explanatory variables, assumed known.

- In the *linear model*, we assume our  $n$  observations (responses) are  $Y_1, \dots, Y_n$  are modelled as

$$Y_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where

- $\beta_1, \dots, \beta_p$  are unknown parameters,  $n > p$
- $x_{i1}, \dots, x_{ip}$  are the values of  $p$  covariates for the  $i$ th response (assumed known)
- $\varepsilon_1, \dots, \varepsilon_n$  are independent (or possibly just uncorrelated) random variables with mean 0 and variance  $\sigma^2$ .

From (??),

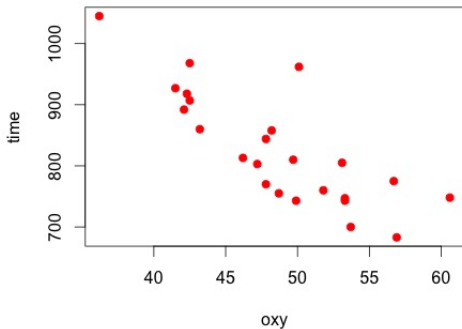
- $\mathbb{E}(Y_i) = \beta_1 x_{i1} + \dots + \beta_p x_{ip}$
- $\text{var}(Y_i) = \text{var}(\varepsilon_i) = \sigma^2$
- $Y_1, \dots, Y_n$  are independent (or uncorrelated).

Note that (??) is linear in the parameters  $\beta_1, \dots, \beta_p$  (there are a wide range of more complex models which are non-linear in the parameters).

### Example 12.1

For each of 24 males, the maximum volume of oxygen uptake in the blood and the time taken to run 2 miles (in minutes) were measured. Interest lies on how the time to run 2 miles depends on the oxygen uptake.

```
oxy=c(42.3,53.1,42.1,50.1,42.5,42.5,47.8,49.9,
      36.2,49.7,41.5,46.2,48.2,43.2,51.8,53.3,
      53.3,47.2,56.9,47.8,48.7,53.7,60.6,56.7)
time=c(918, 805, 892, 962, 968, 907, 770, 743,
       1045, 810, 927, 813, 858, 860, 760, 747,
       743, 803, 683, 844, 755, 700, 748, 775)
plot(oxy, time)
```



- For individual  $i$ , let  $Y_i$  be the time to run 2 miles, and  $x_i$  be the maximum volume of oxygen uptake,  $i = 1, \dots, 24$ .
- A possible model is

$$Y_i = a + bx_i + \varepsilon_i, \quad i = 1, \dots, 24,$$

where  $\varepsilon_i$  are independent random variables with variance  $\sigma^2$ , and  $a$  and  $b$  are constants.

# Matrix formulation

The linear model may be written in matrix form. Let

$$\mathbf{Y}_{n \times 1} = \begin{pmatrix} Y_1 \\ \cdot \\ \cdot \\ Y_n \end{pmatrix}, \quad \mathbf{X}_{n \times p} = \begin{pmatrix} x_{11} & \cdot & \cdot & x_{1p} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ x_{n1} & \cdot & \cdot & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta}_{p \times 1} = \begin{pmatrix} \beta_1 \\ \cdot \\ \cdot \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\varepsilon}_{n \times 1} = \begin{pmatrix} \varepsilon_1 \\ \cdot \\ \cdot \\ \varepsilon_n \end{pmatrix},$$

Then from (??),

$$\begin{aligned} \mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ \mathbb{E}(\boldsymbol{\varepsilon}) &= \mathbf{0} \\ \text{cov}(\mathbf{Y}) &= \sigma^2 \mathbf{I} \end{aligned} \tag{2}$$

We assume throughout that  $X$  has full rank  $p$ .

We also assume the error variance is the same for each observation: this is the *homoscedastic* case (as opposed to *heteroscedastic*).

**Example 12.1 continued**

- Recall  $Y_i = a + bx_i + \varepsilon_i$ ,  $i = 1, \dots, 24$ .
- In matrix form:

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \cdot \\ \cdot \\ Y_{24} \end{pmatrix}, \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \cdot & \cdot \\ \cdot & \cdot \\ 1 & x_{24} \end{pmatrix}, \boldsymbol{\beta} = \begin{pmatrix} a \\ b \end{pmatrix}, \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \cdot \\ \cdot \\ \varepsilon_{24} \end{pmatrix},$$

- Then

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

# Least squares estimation

- In a linear model  $\mathbf{Y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ , the *least squares estimator*  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  minimises

$$\begin{aligned} S(\boldsymbol{\beta}) &= \|\mathbf{Y} - X\boldsymbol{\beta}\|^2 = (\mathbf{Y} - X\boldsymbol{\beta})^T (\mathbf{Y} - X\boldsymbol{\beta}) \\ &= \sum_{i=1}^n (Y_i - \sum_{j=1}^p x_{ij}\beta_j)^2 \end{aligned}$$

- So

$$\left. \frac{\partial S}{\partial \beta_k} \right|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} = 0, \quad k = 1, \dots, p.$$

- So  $-2 \sum_{i=1}^n x_{ik} (Y_i - \sum_{j=1}^p x_{ij}\hat{\beta}_j) = 0$ ,  $k = 1, \dots, p$ .
- i.e.  $\sum_{i=1}^n x_{ik} \sum_{j=1}^p x_{ij}\hat{\beta}_j = \sum_{i=1}^n x_{ik} Y_i$ ,  $k = 1, \dots, p$ .
- In matrix form,

$$X^T X \hat{\boldsymbol{\beta}} = X^T \mathbf{Y} \tag{3}$$

the *least squares equation*.



- Recall we assume  $X$  is of full rank  $p$ .
- This implies

$$\mathbf{t}^T X^T X \mathbf{t} = (X\mathbf{t})^T (X\mathbf{t}) = \|X\mathbf{t}\|^2 > 0$$

for  $\mathbf{t} \neq \mathbf{0}$  in  $\mathbb{R}^p$ .

- i.e.  $X^T X$  is *positive definite*, and hence has an inverse.
- Hence

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{Y} \quad (4)$$

which is linear in the  $Y_i$ 's.

- We also have that

$$\mathbb{E}(\hat{\beta}) = (X^T X)^{-1} X^T \mathbb{E}(\mathbf{Y}) = (X^T X)^{-1} X^T X \beta = \beta$$

so  $\hat{\beta}$  is unbiased for  $\beta$ .

- And

$$\text{cov}(\hat{\beta}) = (X^T X)^{-1} X^T \text{cov}(\mathbf{Y}) X (X^T X)^{-1} = (X^T X)^{-1} \sigma^2 \quad (5)$$

since  $\text{cov}(\mathbf{Y}) = \sigma^2 I$ .

# Simple linear regression using standardised $x$ 's

- The model

$$Y_i = a + bx_i + \varepsilon_i, \quad i = 1, \dots, n,$$

can be reparametrised to

$$Y_i = a' + b(x_i - \bar{x}) + \varepsilon_i, \quad i = 1, \dots, n, \quad (6)$$

where  $\bar{x} = \sum x_i/n$  and  $a' = a + b\bar{x}$ .

- Since  $\sum_i (x_i - \bar{x}) = 0$ , this leads to simplified calculations.

- In matrix form,  $X = \begin{pmatrix} 1 & (x_1 - \bar{x}) \\ \cdot & \cdot \\ 1 & (x_{24} - \bar{x}) \end{pmatrix}$ , so that  $X^T X = \begin{pmatrix} n & 0 \\ 0 & S_{xx} \end{pmatrix}$ ,  
where  $S_{xx} = \sum_i (x_i - \bar{x})^2$ .

- Hence

$$(X^T X)^{-1} = \begin{pmatrix} \frac{1}{n} & 0 \\ 0 & \frac{1}{S_{xx}} \end{pmatrix},$$

so that

$$\hat{\beta} = (X^T X)^{-1} X^T \mathbf{Y} = \begin{pmatrix} \bar{Y} \\ \frac{S_{xY}}{S_{xx}} \end{pmatrix},$$

where  $S_{xY} = \sum_i Y_i (x_i - \bar{x})$ .

- We note that the estimated intercept is  $\hat{a}' = \bar{y}$ , and the estimated gradient  $\hat{b}$  is

$$\begin{aligned}\hat{b} &= \frac{S_{xy}}{S_{xx}} = \frac{\sum_i y_i(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} = \frac{\sum_i (y_i - \bar{y})(x_i - \bar{x})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}} \times \sqrt{\frac{S_{yy}}{S_{xx}}} \\ &= r \times \sqrt{\frac{S_{yy}}{S_{xx}}}\end{aligned}$$

- Thus the estimated gradient is the *Pearson product-moment correlation coefficient*  $r$ , times the ratio of the empirical standard deviations of the  $y$ 's and  $x$ 's.

(Note this estimated gradient is the same whether the  $x$ 's are standardised to have mean 0 or not.)

- From (5),  $\text{cov}(\hat{\beta}) = (X^T X)^{-1} \sigma^2$ , and so

$$\text{var}(\hat{a}') = \text{var}(\bar{Y}) = \frac{\sigma^2}{n}; \quad \text{var}(\hat{b}) = \frac{\sigma^2}{S_{xx}};$$

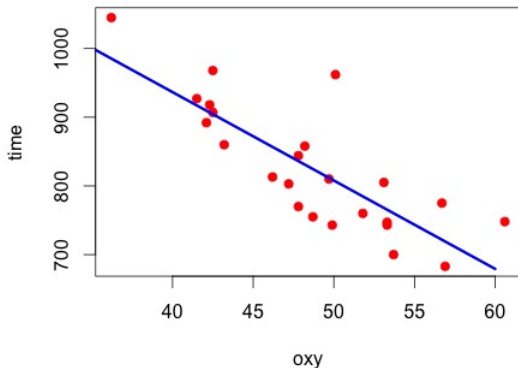
- These estimators are uncorrelated.

All these results are obtained without any explicit distributional assumptions.

**Example 12.1 continued**

$$n = 24, \hat{a}' = \bar{y} = 826.5.$$

$$S_{xx} = 783.5 = 28.0^2, S_{xy} = -10077, S_{yy} = 444^2, r = -0.81, \hat{b} = -12.9.$$



Line goes through  $(\bar{x}, \bar{y})$ .

# 'Gauss Markov' theorem

## Theorem 12.2

In the full rank linear model, let  $\hat{\beta}$  be the least squares estimator of  $\beta$  and let  $\beta^*$  be any other unbiased estimator for  $\beta$  which is linear in the  $Y_i$ 's.

Then  $\text{var}(\mathbf{t}^T \hat{\beta}) \leq \text{var}(\mathbf{t}^T \beta^*)$  for all  $\mathbf{t} \in \mathbb{R}^p$ .

We say that  $\hat{\beta}$  is the Best Linear Unbiased Estimator of  $\beta$  (BLUE).

### Proof:

- Since  $\beta^*$  is linear in the  $Y_i$ 's,  $\beta^* = A\mathbf{Y}$  for some  $A_{p \times n}$ .
- Since  $\beta^*$  is unbiased, we have that  $\beta = \mathbb{E}(\beta^*) = AX\beta$  for all  $\beta \in \mathbb{R}^p$ , and so  $AX = I_p$ .
- Now

$$\begin{aligned} \text{cov}(\beta^*) &= \mathbb{E}(\beta^* - \beta)(\beta^* - \beta)^T) \\ &= \mathbb{E}(AX\beta + A\epsilon - \beta)(AX\beta + A\epsilon - \beta)^T) \\ &= \mathbb{E}(A\epsilon\epsilon^T A^T) \quad \text{since } AX\beta = \beta \\ &= A(\sigma^2 I)A^T = \sigma^2 AA^T \end{aligned}$$

- Now  $\beta^* - \hat{\beta} = (A - (X^T X)^{-1} X^T) \mathbf{Y} = B \mathbf{Y}$ , say.
- And  $BX = AX - (X^T X)^{-1} X^T X = I_p - I_p = 0$ .
- So

$$\begin{aligned} \text{cov}(\beta^*) &= \sigma^2 (B + (X^T X)^{-1} X^T) (B + (X^T X)^{-1} X^T)^T \\ &= \sigma^2 (BB^T + (X^T X)^{-1}) \\ &= \sigma^2 BB^T + \text{cov}(\hat{\beta}) \end{aligned}$$

- So for  $\mathbf{t} \in \mathbb{R}^p$ ,

$$\begin{aligned} \text{var}(\mathbf{t}^T \beta^*) &= \mathbf{t}^T \text{cov}(\beta^*) \mathbf{t} = \mathbf{t}^T \text{cov}(\hat{\beta}) \mathbf{t} + \mathbf{t}^T BB^T \mathbf{t} \sigma^2 \\ &= \text{var}(\mathbf{t}^T \hat{\beta}) + \sigma^2 \|B^T \mathbf{t}\|^2 \\ &\geq \text{var}(\mathbf{t}^T \hat{\beta}). \end{aligned}$$

- Taking  $\mathbf{t} = (0, \dots, 1, 0, \dots, 0)^T$  with a 1 in the  $i$ th position, gives

$$\text{var}(\hat{\beta}_i) \leq \text{var}(\beta_i^*).$$

□

# Fitted values and residuals

## Definition 12.3

- $\hat{\mathbf{Y}} = X\hat{\beta}$  is the vector of *fitted values*.
  - $\mathbf{R} = \mathbf{Y} - \hat{\mathbf{Y}}$  is the vector of *residuals*.
  - The *residual sum of squares* is  $\text{RSS} = \|\mathbf{R}\|^2 = \mathbf{R}^T \mathbf{R} = (\mathbf{Y} - X\hat{\beta})^T (\mathbf{Y} - X\hat{\beta})$
- 
- Note  $X^T \mathbf{R} = X^T (\mathbf{Y} - \hat{\mathbf{Y}}) = X^T \mathbf{Y} - X^T X \hat{\beta} = 0$  by (??).
  - So  $\mathbf{R}$  is orthogonal to the column space of  $X$ .
  - Write  $\hat{\mathbf{Y}} = X\hat{\beta} = X(X^T X)^{-1} X^T \mathbf{Y} = P\mathbf{Y}$ , where  $P = X(X^T X)^{-1} X^T$ .
  - $P$  represents an orthogonal projection of  $\mathbb{R}^n$  onto the space spanned by columns of  $X$ . We have  $P^2 = P$  ( $P$  is idempotent) and  $P^T = P$  (symmetric).

