# 4 Linear regression

## 4.1 Simple LR    See notes

## 4.2 Multiple LR

Setup: Features $X_1, \ldots, X_p$, continuous response $Y$

MLR model is:

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon \quad \longleftarrow \text{random} \tag{3}$$

We say $Y$ is being <u>regressed</u> on $X_1, \ldots, X_p$. Here, $\beta_0, \ldots, \beta_p$ are <u>regression parameters</u>. This is the usual model

$$Y = f(\underline{x}) + \varepsilon$$

with a linear specification for $f$: $f(\underline{x}) = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$

$f$ is the <u>population regression line</u> & $\varepsilon$ is the <u>residual</u>.

$\beta_0$ = Avg value of $y$ when $X_1 = \ldots = X_p = 0$

$\beta_k$ = Avg change in $Y$ for a unit increase in $X_k$ with all other features fixed

Given a set of data $Y_1, \ldots, Y_n$ with $y_i$ having corresponding features $X_{i1}, X_{i2}, \ldots, X_{ip}$, we have $n$
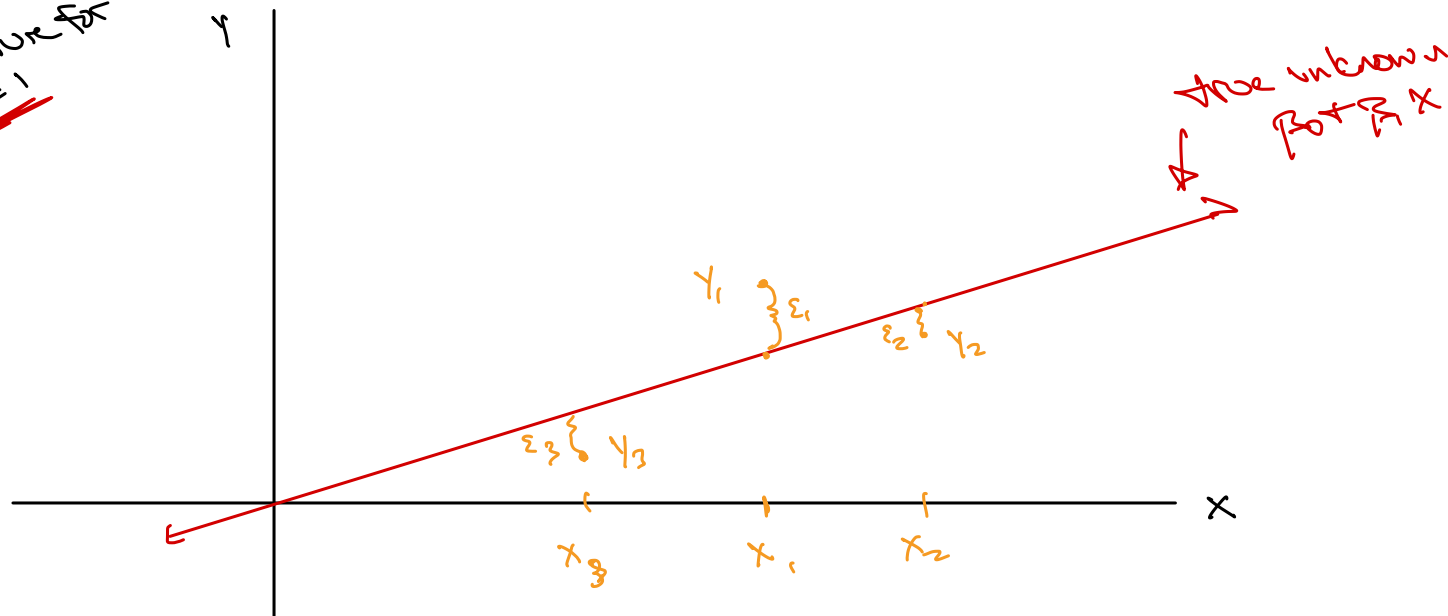
<u>observation equations</u>

$$Y_1 = \beta_0 + \beta_1 X_{11} + \cdots + \beta_p X_{1p} + \varepsilon_1$$

$$Y_2 = \beta_0 + \beta_1 X_{21} + \cdots + \beta_p X_{2p} + \varepsilon_2$$

$$\vdots$$

$$Y_n = \beta_0 + \beta_1 X_{n1} + \cdots + \beta_p X_{np} + \varepsilon_n$$

Picture for
$P = 1$

$X_1$

true unknown
$\beta_0 + \beta_1 X$

$Y_1$   $\varepsilon_1$

$\varepsilon_2$   $Y_2$

$\varepsilon_3$   $Y_3$

$X_3$   $X_1$   $X_2$

Annoying to write

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip} + \varepsilon_i \quad \text{for } i = 1, \ldots, n$$

$$\underset{n \times 1}{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} \qquad \underset{(p+1) \times 1}{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \qquad \underset{n \times 1}{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\underset{n \times (p+1)}{X} = \begin{matrix} \beta_0 & \beta_1 & \beta_2 & \cdots & \beta_p \end{matrix}\begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}$$

$$\underset{n \times 1}{Y} = \underset{n \times (p+1)}{X} \underset{(p+1) \times 1}{\beta} + \underset{n \times 1}{\varepsilon} \qquad (4)$$

# 4.2 Assumptions

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_p x_{ip} + \varepsilon_i \quad , i = 1, \ldots, n \quad (4)$$

What do we assume about (4)?

**A1** · (4) holds · $\{\varepsilon_i\}$ iid $N(0, \sigma^2)$

$\boxed{A2}$ · (4) holds · $E\varepsilon_i = 0$ $\forall i$ · $\text{Var } \varepsilon_i = \sigma^2$ $\forall i$

<span style="color:red">(homoskedastic)</span>

· $\{\varepsilon_i\}$ are iid

$[$ relaxing normality $]$

$\boxed{A3}$ · (4) holds · $E\varepsilon_i = 0$ $\forall i$ · $\text{Var } \varepsilon_i = \sigma^2$ $\forall i$

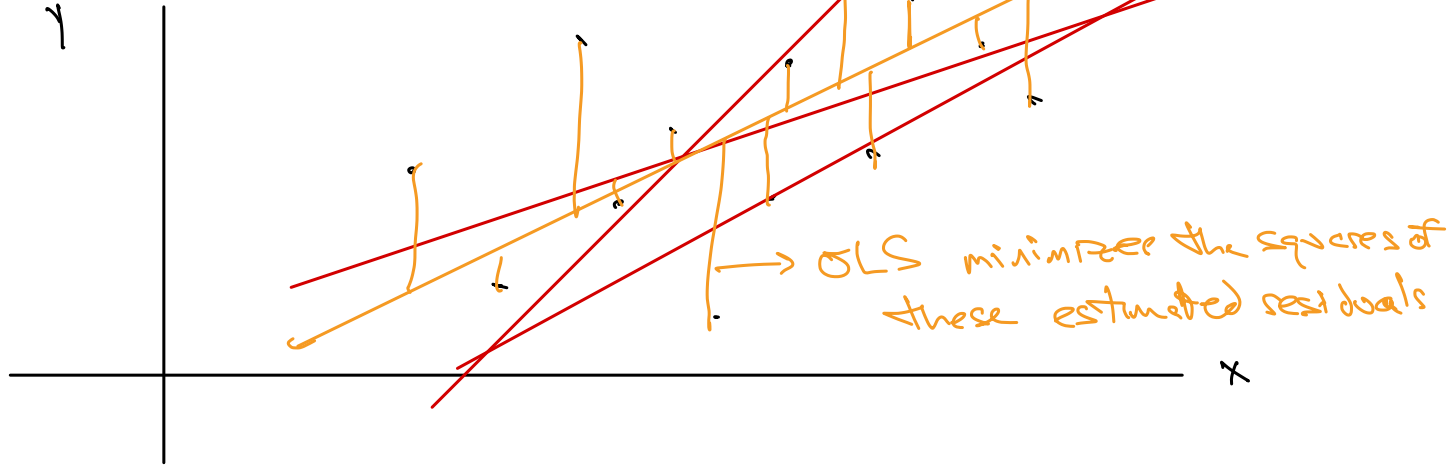· $\{\varepsilon_i\}$ pairwise uncorrelated

$[$ relaxing iid $]$

---

We will always assume $\boxed{A1} - \boxed{A3}$ , in which case

$$E \underline{\varepsilon} = \underline{0} \qquad \text{Var } \underline{\varepsilon} = \sigma^2 I$$

so $\Rightarrow$ $E\underline{Y} = E(X\underline{\beta} + \underline{\varepsilon}) = E(X\underline{\beta}) + E\underline{\varepsilon} = X\underline{\beta}$

$\text{Var } \underline{Y} = \text{Cov}(\underline{Y}, \underline{Y}) = \text{Cov}(X\underline{\beta} + \underline{\varepsilon}, X\underline{\beta} + \underline{\varepsilon}) = \text{Cov}(\underline{\varepsilon}, \underline{\varepsilon}) = \text{Var } \underline{\varepsilon} = \sigma^2 I$

$y$

→ OLS minimizes the squares of these estimated residuals

$x$

The OLS estimator for $\beta$ minimizes residual sum of squares:

$$RSS(\beta) = \sum_{i=1}^{n} \left( y_i - (\beta_0 + \beta_1 x_{i,1} + \cdots + \beta_p x_{i,p}) \right)^2$$

$$= (y - X\beta)^\top (y - X\beta)$$

The OLS estimator is:

$$\hat{\beta} = \hat{\beta}_{OLS} = (X^TX)^{-1}X^T y$$

**Note!**

$$E\,\hat{\beta} = E\,(X^TX)^{-1}X^T y = (X^TX)^{-1}X^T\,Ey$$

$$= (X^TX)^{-1}X^TX\,\beta = I\beta = \beta \qquad \Rightarrow \hat{\beta} \text{ is unbiased}$$

$$Var\,\hat{\beta} = Cov\left(\hat{\beta},\hat{\beta}\right) = Cov\left((X^TX)^{-1}X^T y,\, (X^TX)^{-1}X^T y\right)_{\,T}$$

$$= (X^TX)^{-1}X^T\,Cov(y,y)\,X^{TT}(X^TX)^{-T}$$

$$= (X^TX)^{-1}X^T\underbrace{(Var\,y)}_{\sigma^2 I}X(X^TX)^{-1}$$

$$= \sigma^2 (X^TX)^{-1} X^TX (X^TX)^{-1}$$

$$\boxed{\text{Var } \hat{\beta} = \sigma^2 (X^TX)^{-1}} \quad \rightarrow \quad (p+1) \times (p+1)$$

$(p+1) \times 1$

$(p+1) \times n$

$n \times (p+1)$

Standard errors for $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ are the square roots of the diagonal of

$$\widehat{\text{Var } \hat{\beta}} = \hat{\sigma}^2 (X^TX)^{-1}$$

where, for example,

$$\hat{\sigma}^2 = \frac{1}{n-(p+1)} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$$

Approx 95% CI for $\beta_j$ is

$$\hat{\beta}_j \pm 2 \, SE(\hat{\beta}_j)$$