

M3L1: linear regression

Simple linear regression

one x

Model: $Y = f(x) + \epsilon = \beta_0 + \beta_1 \cdot x + \epsilon$

Annotations: β_0 is intercept, β_1 is slope, ϵ is error term. x is observed.

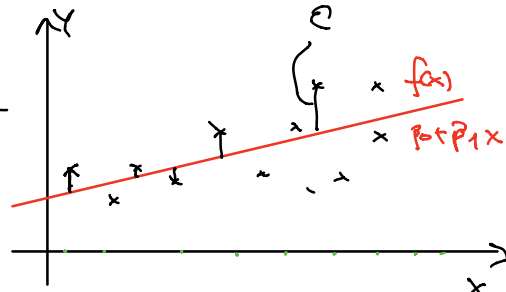
Additional assumption:

ϵ has $E(\epsilon) = 0$

$Var(\epsilon) = \sigma^2$

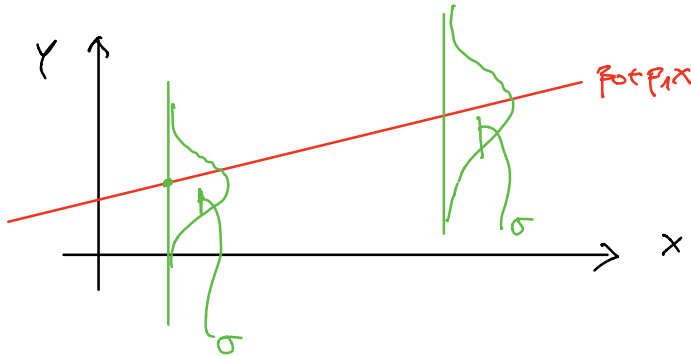
not dependent on x

homoscedastic error



(if $Var(\epsilon)$ varies with x \Rightarrow heteroscedastic error)

And often $\epsilon \sim N(0, \sigma^2)$



And we also assume that the pairs $(x_i, y_i) \quad i=1, \dots, n$ are independent

PARAMETER ESTIMATION

For a given dataset $(x_i, y_i), i=1, \dots, n$ independent pairs.

We don't know β_0, β_1 and σ^2 , and need to find estimators

Annotations: β_0, β_1 are parameters unknown; σ^2 is observed.

$\hat{\beta}_0, \hat{\beta}_1$ and $\hat{\sigma}^2$

- least squares
- maximum likelihood

restricted maximum likelihood

Let $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$. We find $\hat{\beta}_0$ and $\hat{\beta}_1$ by minimizing

$$RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (Y_i = \beta_0 + \beta_1 x_i + \epsilon_i)$$

residual
sum of squares

e_i
residual
predicted value for the error ϵ_i
 $e_i = \hat{\epsilon}_i$ RV

$$\left. \begin{aligned} \frac{\partial RSS}{\partial \beta_0} &= 0 \\ \frac{\partial RSS}{\partial \beta_1} &= 0 \end{aligned} \right\} \begin{aligned} \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X} \\ \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

We also need $\hat{\sigma}^2$ ← the variance of ϵ 's
 $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$, $\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$

Remember that $\epsilon \sim N(0, \sigma^2)$; $\text{Var}(\epsilon) = E(\epsilon^2) - \underbrace{E(\epsilon)}_0^2$

Since the $\hat{\epsilon}_i$'s are predictions of the ϵ_i 's, we use them.

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n \hat{\epsilon}_i^2}{n-2} = \frac{RSS}{n-2}$$

$$RSE = \hat{\sigma}$$

#param. estimated
($\hat{\beta}_0, \hat{\beta}_1$)

Distribution of parameter estimators

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x}) Y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}$$

$$Y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\text{constants}} + \epsilon_i \quad \begin{matrix} \swarrow N(0, \sigma^2) \\ \Rightarrow Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2) \end{matrix}$$

$$E(Y_i) = \beta_0 + \beta_1 x_i + \underbrace{E(\epsilon_i)}_0$$

$$\hat{\beta}_j \sim N(\beta_j, \underbrace{\text{Var}(\hat{\beta}_j)}_{c_{jj}\sigma^2})$$

$j=0,1$

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{c_{jj}}\sigma} \sim N(0,1)$$

+ can be shown

$$\frac{\hat{\sigma}^2(n-2)}{\sigma^2} \sim \chi^2_{n-2}$$

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{c_{jj}}\hat{\sigma}} \sim t_{n-2}$$

← $n - \overset{\text{regr.}}{\text{\#param estimated}}$

$(1-\alpha) \cdot 100\%$
 CI : $\left[\hat{\beta}_j \pm t_{\frac{\alpha}{2}, n-2} \cdot \sqrt{c_{jj}} \cdot \hat{\sigma} \right]$

confidence intervals

Hypothesis test: $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$

	reject H_0	not reject H_0
H_0 true	type I error	correct
H_0 false	correct	type II error

crime of justice

guilty criminal go free

$P(\text{type I error}) \leq \alpha$
 0.05

p-value = $P(T_0 \geq |t_0| \mid H_0 \text{ true})$
 reject H_0 when p-value $\leq \alpha$. 3

How good is the regression

$$TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \text{total variability}$$

↑
total sums of squares

$$RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \text{not explained by regression}$$

regression line

also explained by regression

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS} \in [0, 1]$$

High R^2 is good.

Multiple linear regression (MLR)

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p \cdot X_{ip} + \epsilon_i$$

$$\begin{cases} Y_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \dots + \epsilon_1 \\ Y_2 = \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \dots + \epsilon_2 \\ \vdots \\ Y_n = \beta_0 + \beta_1 X_{n1} + \beta_2 X_{n2} + \dots + \epsilon_n \end{cases}$$

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\begin{array}{ccccccc} Y & = & X & \cdot & \beta & + & \epsilon \\ n \times 1 & & n \times (p+1) & & (p+1) \times 1 & & n \times 1 \\ \uparrow & & \uparrow & & \uparrow & & \uparrow \\ \text{response} & & \text{design} & & \text{vector of} & & \text{vector of} \\ \text{vector} & & \text{matrix} & & \text{regression} & & \text{errors} \\ & & & & \text{param} & & \end{array}$$

Combine independent pairs (x_i, y_i) and $E(\epsilon_i) = 0, \text{Var}(\epsilon_i) = \sigma^2$

$$\Rightarrow \epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix} \sim N_n(0, \sigma^2 I)$$

Multivariate normal
dim n

$$\begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & & \\ & & \ddots & \\ 0 & & & \sigma^2 \end{bmatrix}$$

Homework: Distribution of Y?