

## Linear Regression Models

- scatterplot  $Y$  vs  $x$
- transformation ( $\log$ ,  $\sqrt{\cdot}$ , ...)

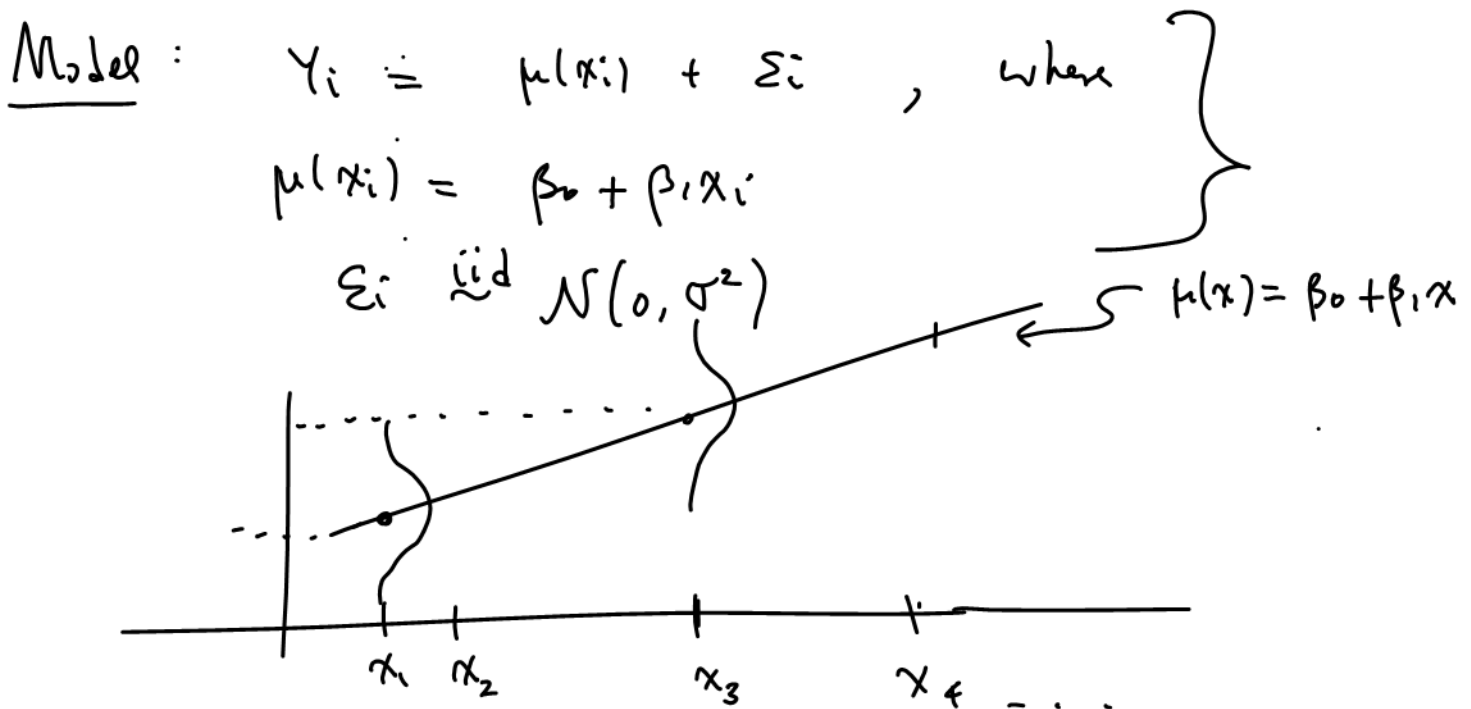
### Random Data

$$\{(x_i, y_i), i=1, \dots, n\}$$

### Observed data

$$\{(x_i, y_i), i=1, \dots, n\}$$

treat  $\{x_i, i=1, \dots, n\}$  fixed & known.



- $y_i | x_i \sim \mathcal{N}(\mu(x_i) = \beta_0 + \beta_1 x_i, \sigma^2)$
- $\text{Var}(y_i | x_i) = \sigma^2 \quad \forall i$  constant across all values of  $x$

- For a fixed covariate  $x^*$ , 99% of all values of the outcome variable  $Y$  fall in the interval  $[\mu(x^*) - 3\sigma, \mu(x^*) + 3\sigma]$

- Mean function  $\mu(x_i)$

$$E(Y_i | x_i) = \mu(x_i) = \beta_0 + \beta_1 x_i \quad \leftarrow$$

Interpretation of  $\beta_0$

$$\text{when } x_i = 0 \Rightarrow \mu(x_i) = \beta_0 \quad \leftarrow$$

$$\therefore \beta_0 = E(Y | x=0)$$

Rem:  $z_i = x_i - \bar{x}$

$$E(Y_i | z_i) = \beta_0 + \beta_1 z_i$$

$$\beta_0 = E(Y_i | z_i = 0) = E(Y_i | x_i = \bar{x})$$

Interpretation of  $\beta_1$

$$E(Y_i | x_i) = \mu(x_i) = \beta_0 + \beta_1 x_i$$

$$E(Y_i | x_i = a+1) = \beta_0 + \beta_1 (a+1)$$

$$E(Y_i | x_i = a) = \beta_0 + \beta_1 a$$

The least squares estimator for  $\beta_0, \beta_1$ :

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$= \sum w_i y_i \quad \text{where } w_i = \frac{x_i - \bar{x}}{\sum (x_j - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

The estimate for  $\sigma^2$  is:

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2}{n - 2}$$

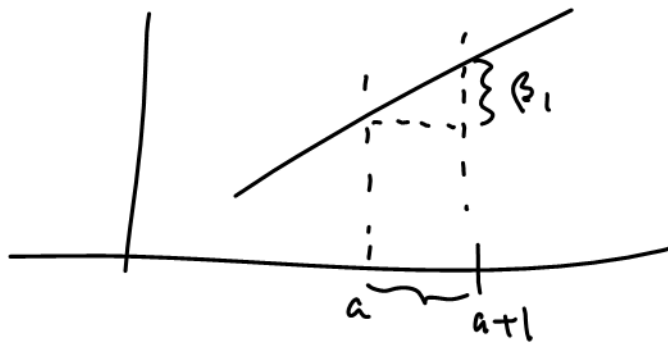
$\nwarrow$  # parameters in  $\mu(x_i)$

---

Inference on  $\beta_1$

Test  $H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$

$$\Rightarrow E(y_i | x_i = a+1) - E(y_i | x_i = a) \\ = \mu(a+1) - \mu(a) = \beta_1$$



$\beta_1$  = change in the mean value of the dist<sup>n</sup> of the outcome variable for every unit increase in the independent variable

Construct a procedure for testing  $H_0: \beta_1 = 0$  vs  $H_1: \beta_1 \neq 0$

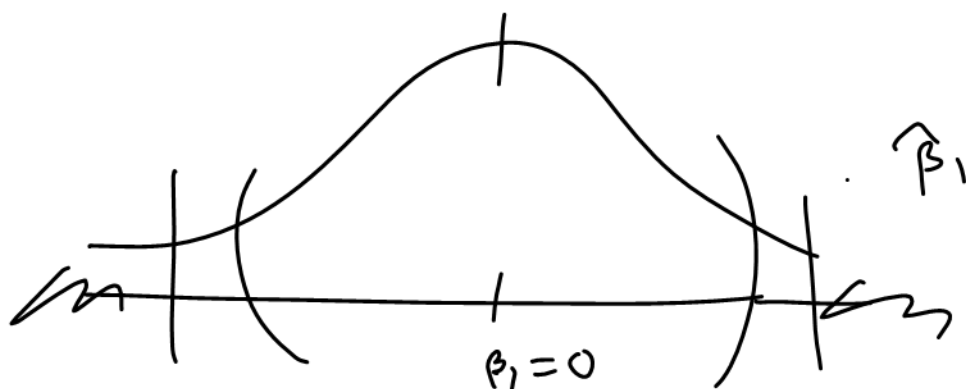
Note: When  $\beta_1 = 0 \Rightarrow E(Y_i | x_i) = \beta_0$  constant over all  $x$

We use  $\hat{\beta}_1$ :

$$\hat{\beta}_1 = \sum w_i Y_i$$

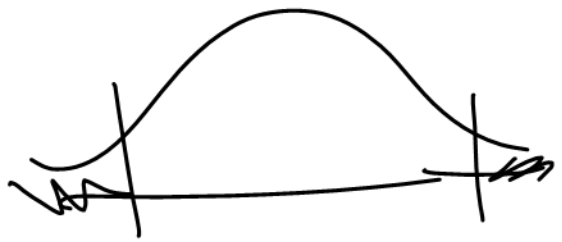
$$\sim N\left(E\hat{\beta}_1 = \beta_1, \text{Var}\hat{\beta}_1 = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right)$$

Population:  $\{(x_i, Y_i), i = 1, \dots, \infty\}$



Idea: Reject  $H_0$  if  $|\hat{\beta}_1|$  is "large".

We need a reference distn:



$$|\hat{\beta}_1| > 2 \cdot \text{SD}(\hat{\beta}_1) \\ > 2 \cdot \sqrt{\frac{\sigma^2}{\sum (x_i - \bar{x})^2}}$$

Problem: the threshold  $2 \sqrt{\frac{\sigma^2}{\sum (x_i - \bar{x})^2}}$  is not known since  $\sigma^2$  is not known.

Results:  $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{\sum (x_i - \bar{x})^2})$

$$\Rightarrow \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{\sigma^2}{\sum (x_i - \bar{x})^2}}} \sim N(0, 1) \quad (A)$$

$$\hat{\sigma}^2 = \frac{\sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2}{n-2} \quad (B)$$

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} = \frac{\sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2}{\sigma^2} \sim \chi^2 \quad (df = n-2)$$

$$\frac{y_i - \mu(x_i)}{\sigma} \sim N(0, 1)$$

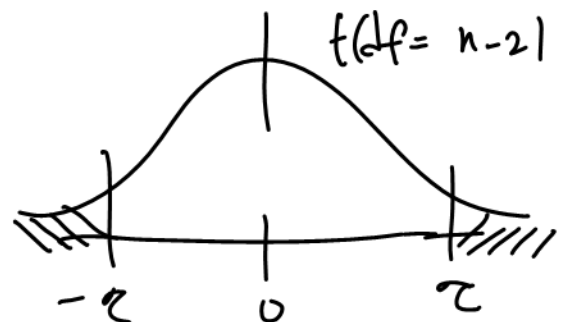
$\hat{\beta}_1, \hat{\sigma}^2$  statistically independent (C)

$$\Rightarrow \frac{\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 / \sum (x_i - \bar{x})^2}}}{\sqrt{\frac{(n-2)\hat{\sigma}^2}{\hat{\sigma}^2} / (n-2)}} \sim t(df = n-2)$$

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2 / \sum (x_i - \bar{x})^2}} \sim t(df = n-2)$$

$$T = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2 / \sum (x_i - \bar{x})^2}} \sim t(df = n-2)$$

$$T = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2 / \sum (x_i - \bar{x})^2}}$$



Suppose we let  $P(\text{Type I error}) = \alpha = .05$

$$\Rightarrow c = t(0.975, df = n-2)$$

Decision Rule: Reject  $H_0$  at level  $\alpha$  if

$$|T| > \tau = t(.975, df = n-2)$$

$$\Leftrightarrow T > \tau \text{ or } T < -\tau$$

Now we have the observed data:

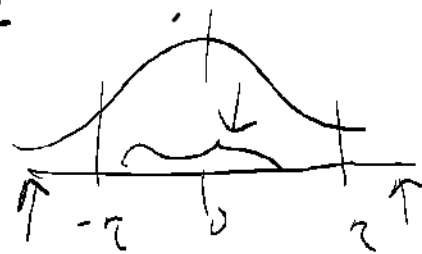
$$\{(x_i, y_i), i = 1, \dots, n\}$$

$$\Rightarrow \begin{aligned} \bar{x}, \bar{y} &= \frac{1}{n} \sum (x_i, y_i) \\ S_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) \\ S_{xx} &= \sum (x_i - \bar{x})^2 \end{aligned}$$

$$\hat{\beta}_{1, obs} = \frac{S_{xy}}{S_{xx}}, \quad \hat{\beta}_{0, obs} = \bar{y} - \hat{\beta}_{1, obs} \bar{x}$$

$$\hat{\sigma}_{obs}^2 = \frac{\sum (y_i - (\hat{\beta}_{0, obs} + \hat{\beta}_{1, obs} x_i))^2}{n-2}$$

$$\Rightarrow T_{obs} = \frac{\hat{\beta}_{1, obs}}{\sqrt{\hat{\sigma}_{obs}^2 / \sum (x_i - \bar{x})^2}}$$





So far:

- Model
- Estimators for the parameters
- Sampling behavior of  $\hat{\beta}_1$

Next:

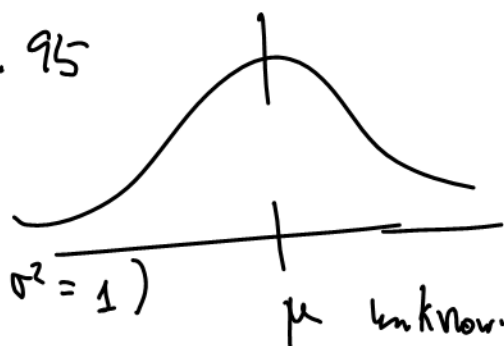
Confidence intervals for  $\beta_1$ ,  $\mu(x^*)$

Let  $\theta$  be an unknown parameter.

Let  $I = [L, U]$  be a random interval estimator for  $\theta$ .

The random interval  $I$  is a 95% <sup>confidence</sup> interval estimator for  $\theta$  if:

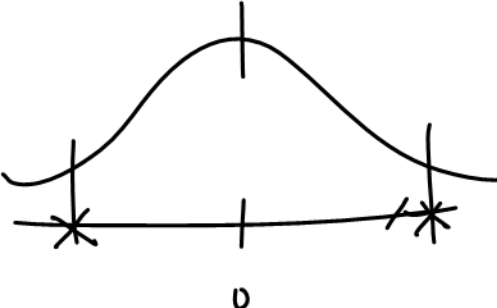
$$P(\theta \in I) = 0.95$$



Let  $Y_1, \dots, Y_n$  i.i.d  $N(\mu, \sigma^2 = 1)$

$\bar{Y} = \frac{1}{n} \sum Y_i$  is the LSE for  $\mu$   $\sigma^2 = 1$   
MLE for  $\mu$

$$\bar{Y} \sim N(\mu, 1/n)$$

$$\Rightarrow \left( \frac{\bar{Y} - \mu}{\sqrt{1/n}} \right) \sim N(0, 1)$$


$$\Rightarrow P\left(-1.96 < \frac{\bar{Y} - \mu}{\sqrt{1/n}} < +1.96\right) = 0.95$$

⋮

$$P(L(\bar{Y}) < \mu < U(\bar{Y})) = 0.95$$

$$P(\mu \in [L(\bar{Y}), U(\bar{Y})]) = 0.95$$

$$\Rightarrow P\left(\bar{Y} - 1.96\left(\frac{1}{\sqrt{n}}\right) < \mu < \bar{Y} + 1.96\left(\frac{1}{\sqrt{n}}\right)\right) = 0.95$$

$$\Rightarrow P\left(\mu \in \left[\bar{Y} - 1.96/\sqrt{n}, \bar{Y} + 1.96/\sqrt{n}\right]\right) = 0.95$$

A 95% CI estimator for  $\mu$  is:

$$\left[ \bar{Y} - 1.96/\sqrt{n}, \bar{Y} + 1.96/\sqrt{n} \right]$$

Observed data  $\Rightarrow$  95% CI estimate  $\left[ \bar{y} - 1.96/\sqrt{n}, \bar{y} + 1.96/\sqrt{n} \right]$ .