

# Linear Regression

## Simple Linear Regression Model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$$

Notation:

$y_i$  =  $i^{\text{th}}$  observed value of the response variable

RANDOM

$x_i$  =  $i^{\text{th}}$  " " of the predictor variable

FIXED

$n$  = sample size

FIXED

$\beta_0$  = unknown intercept parameter

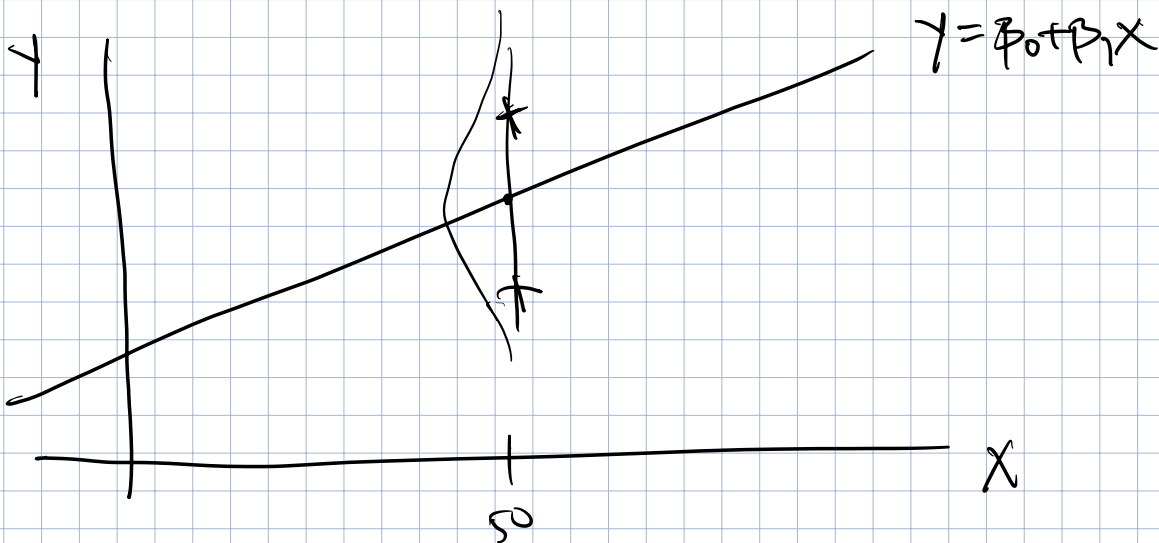
FIXED

$\beta_1$  = unknown slope parameter

FIXED

$\varepsilon_i$  = Random error terms

RANDOM



Observations:

① "simple": one predictor ( $x$ ) variable

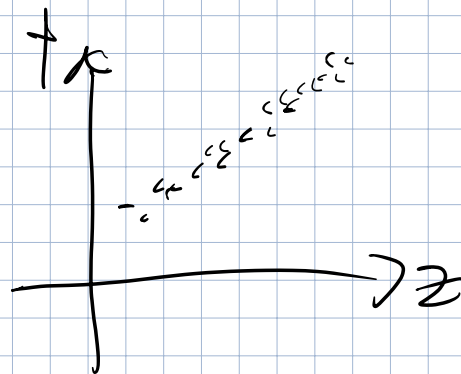
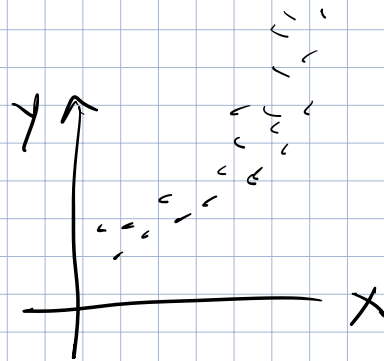
② "linear":  $y_i$  is a linear f'n of  $\beta_0$  &  $\beta_1$

EX:

①  $y_i = \beta_0 + \beta_1 x_i^2 + \varepsilon_i$

$\downarrow$   $z_i = x_i^2$

$$y_i = \beta_0 + \beta_1 z_i + \varepsilon_i$$



②  $y_i = \beta_0 + \beta_1 \log x_i + \varepsilon_i$

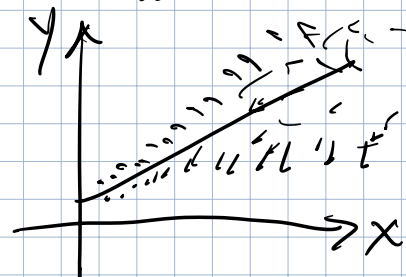
③  $y_i = \frac{\beta_0}{\beta_1 + x_i} + \varepsilon_i$

# Model Assumptions

- $x_i$ 's are fixed

↳ (if not fixed, need random effects model.)

- $\varepsilon_i$ 's must satisfy:



- BASIC
- i.  $E(\varepsilon_i) = 0$
  - ii.  $\text{Var}(\varepsilon_i) = \sigma^2$  (constant in  $X$ )
  - iii.  $\text{Cor}(\varepsilon_i, \varepsilon_j) = 0$  for  $i \neq j$   
↳ errors are uncorrelated

- CLASSICAL
- iv.  $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$   
↳ useful for inference

## 3. Regression Function

The regression function is

$$g(x) = E(Y | X=x)$$

↳ if simple & linear assumptions are made

$$= \beta_0 + \beta_1 x$$

The goal of estimating  $g(x)$  can be reduced to

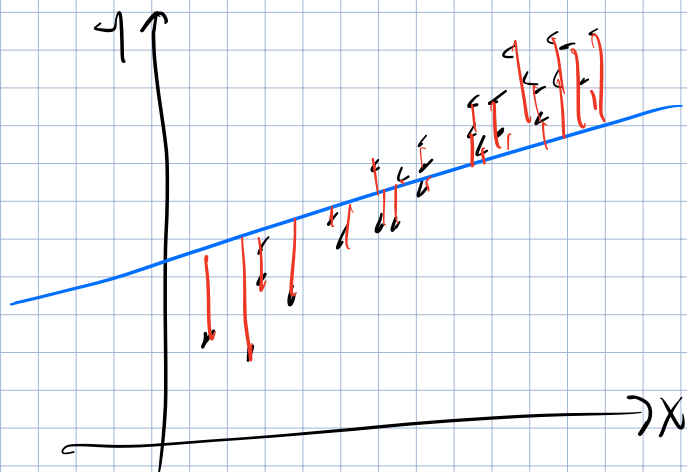
estimating  $\beta_0$  &  $\beta_1$ .

How to estimate  $\beta_0$  &  $\beta_1$ ?

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$\epsilon_i = y_i - \beta_0 - \beta_1 x_i$$

The Least Squares Principle is one way to do this.



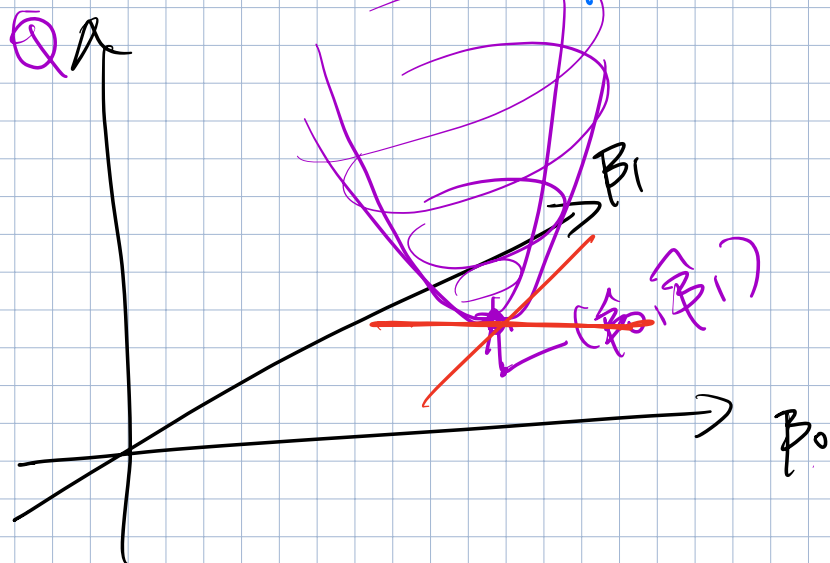
Find the line  $(\beta_0, \beta_1)$  that minimize

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n \epsilon_i^2$$

$$= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

How can I minimize

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$



$$\textcircled{1} \frac{\partial Q}{\partial \beta_0} \stackrel{!}{=} 0$$

$$\textcircled{2} \frac{\partial Q}{\partial \beta_1} \stackrel{!}{=} 0$$

$$\begin{aligned}
 \textcircled{1} \quad \frac{\partial Q}{\partial \beta_0} &= \frac{\partial}{\partial \beta_0} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \\
 &= \sum_{i=1}^n \frac{\partial}{\partial \beta_0} (y_i - \beta_0 - \beta_1 x_i)^2 \\
 &= \sum_{i=1}^n 2 (y_i - \beta_0 - \beta_1 x_i) (0 - 1 - 0)
 \end{aligned}$$

$$= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \stackrel{!}{=} 0$$

what does this imply?

$$\begin{aligned}
 \textcircled{2} \quad \frac{\partial Q}{\partial \beta_1} &= \frac{\partial}{\partial \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \\
 &= \sum_{i=1}^n \frac{\partial}{\partial \beta_1} (y_i - \beta_0 - \beta_1 x_i)^2 \\
 &= \sum_{i=1}^n 2 (y_i - \beta_0 - \beta_1 x_i) (0 - 0 - x_i) \\
 &= -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i \stackrel{!}{=} 0
 \end{aligned}$$

①

$$-2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) \stackrel{!}{=} 0$$

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum_i y_i - \sum_i \beta_0 - \sum_i \beta_1 x_i = 0$$

$$n\bar{y} - n\beta_0 - \beta_1 n\bar{x} = 0$$

$$n\bar{y} = n\beta_0 + n\beta_1 \bar{x}$$

$$\bar{y} = \beta_0 + \beta_1 \bar{x}$$

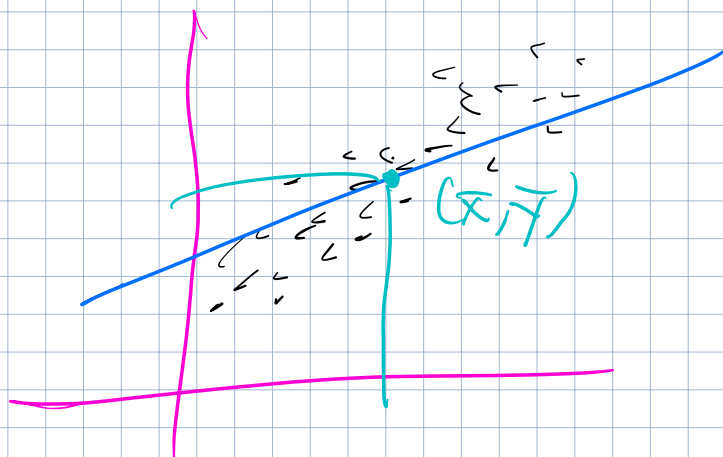
↙

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$$\bar{y} = \frac{\sum_i y_i}{n}$$

$$n\bar{y} = \sum_i y_i$$

↑



②

$$-2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i \stackrel{!}{=} 0$$

$$\sum_{i=1}^n (x_i y_i - \beta_0 x_i - \beta_1 x_i^2) = 0$$

$$\sum_{i=1}^n x_i y_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0$$

$$\sum_i x_i y_i - (\bar{y} - \beta_1 \bar{x}) (\sum_i x_i) - \beta_1 \sum_i x_i^2 = 0$$

$\underbrace{\sum_i x_i}_{n\bar{x}}$

$$\sum_i x_i y_i - (n \bar{x} \bar{y} - n \beta_1 \bar{x}^2) - \beta_1 \sum_i x_i^2 = 0$$

$$\sum_i x_i y_i - n \bar{x} \bar{y} + n \beta_1 \bar{x}^2 - \beta_1 \sum_i x_i^2 = 0$$

$$\begin{aligned} \sum_i x_i y_i - n \bar{x} \bar{y} &= \beta_1 \sum_i x_i^2 - \beta_1 n \bar{x}^2 \\ &= \beta_1 (\sum_i x_i^2 - n \bar{x}^2) \end{aligned}$$

$$\begin{aligned} \hookrightarrow \hat{\beta}_1 &= \frac{\sum_i x_i y_i - n \bar{x} \bar{y}}{\sum_i x_i^2 - n \bar{x}^2} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \end{aligned}$$

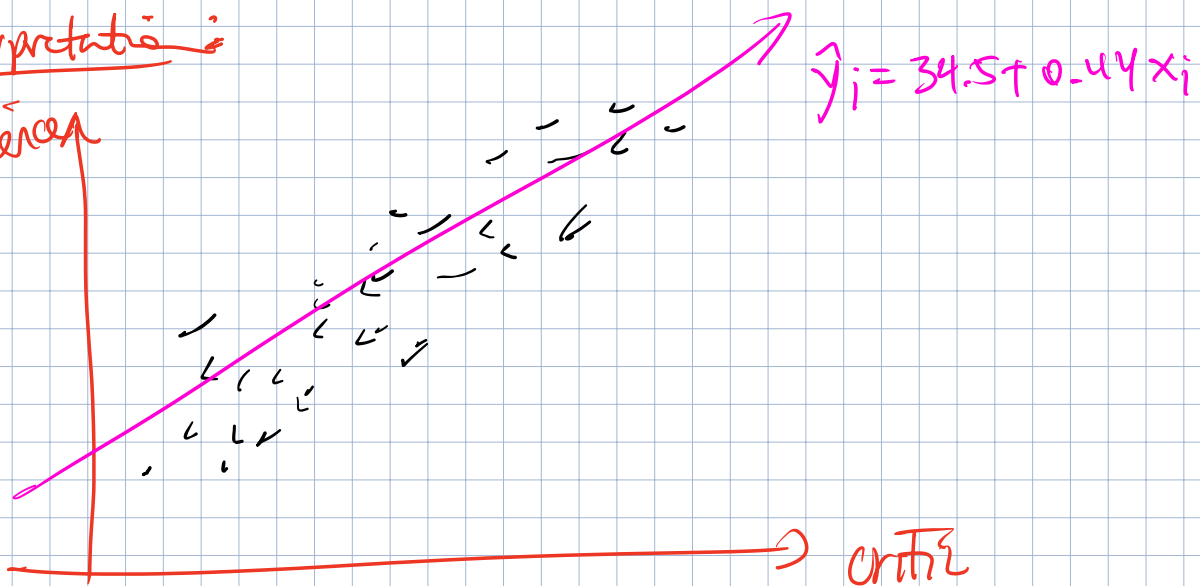
Least Squares Estimates of  $\beta_0$  &  $\beta_1$

~~MM~~

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} = \frac{(n-1) \hat{\text{Cov}}(X, Y)}{(n-1) \hat{\text{Var}}(X)} \\ &= \frac{\frac{1}{n-1} \sum_i (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n-1} \sum_i (x_i - \bar{x})^2} = \frac{\hat{\text{Cov}}(X, Y)}{\hat{\text{Var}}(X)} \end{aligned}$$

Interpretation:

audience



Slope: Interpret 0.44 in the movie ex

"For every additional point a critic gave a movie, we expect the audience rating to increase by 0.44 points, on average."

Intercept: 34.5

We expect on average a movie rated as a 0 by a critic to have an audience rating of 34.5.



