

STA302 - Lecture 1

①

* The statistical model for simple linear regression is:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Y = dependent variable, it's a r.v

X = independent variable

β_0, β_1 = parameters of model, they are unknown and we will use observed data to estimate them.

ε = random error/nois, variation in measures that we can't account for.

* linear in β 's

Simple - only one X ; i.e only one Predictor.

* More formally

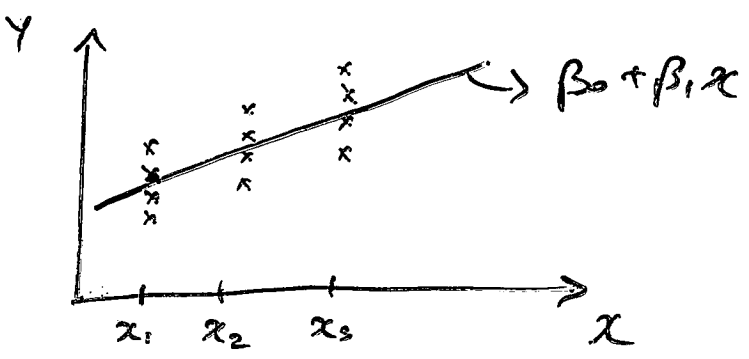
$E(Y|X) = \beta_0 + \beta_1 X$ if X is r.v

$E(Y|X=x) = \beta_0 + \beta_1 x$ if X is fixed.

* Can think about the different values of X as defining different subpopulations one for each possible value of X . Each subpopulation consist of all individuals in the population having the same x value.

* Example:

$Y = \text{weight}$, $X = \text{height}$



(x_i, y_i)
 $i = 1, \dots, 30$

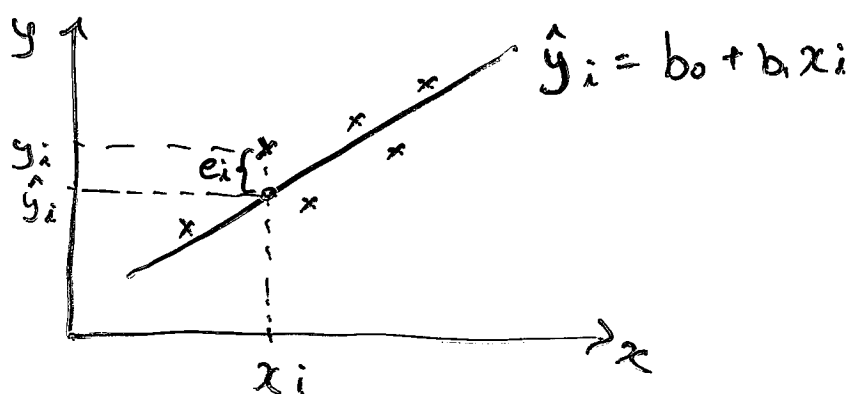
SLR Model - Continuation

3

* Data $(x_i, y_i) \quad i=1, \dots, n$

* Model: $Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad i=1, \dots, n$

* Graphical representation of data is a simple scatterplot.



* Want to estimate β_0 and β_1 based on observed data. Estimated are denoted by b_0, b_1 or $\hat{\beta}_0$ and $\hat{\beta}_1$.

* For each observed value x_i of X the fitted value of Y is $\hat{y}_i = b_0 + b_1 x_i$

Least-Square Estimates

5

* Why vertical distances?

- Want to predict Y from X and so we want \hat{y}_i (the predicted value) to be as close as possible to y_i (observed value)
- If we minimize the horizontal distances, we will get a different answer for b_0 and b_1 .
- Regression is not symmetric - it matters which variable is dependent and which is independent.

* Why minimizing squared deviations?

Answer: Since $\sum_{i=1}^n e_i = 0$ (proof later...)

Derivation of LS estimates

* Residuals, vertical distances:

$$e_i = y_i - \hat{y}_i = y_i - (b_0 + b_1 x_i)$$

\hat{y}_i = the fitted line.

$$* \text{RSS} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

want to find b_0 and b_1 that minimized RSS.

$$* \text{So, set } \frac{\partial \text{RSS}}{\partial b_0} = 0 \text{ and } \frac{\partial \text{RSS}}{\partial b_1} = 0$$

$$\Rightarrow \frac{\partial \text{RSS}}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i)$$

$$\frac{\partial \text{RSS}}{\partial b_1} = -2 \sum_{i=1}^n x_i (y_i - b_0 - b_1 x_i)$$

* So b_0 and b_1 must satisfy:

$$\textcircled{1} \sum_{i=1}^n y_i = n b_0 + b_1 \sum x_i$$

$$\textcircled{2} \sum_{i=1}^n x_i y_i = b_0 \sum x_i + b_1 \sum x_i^2$$

} "Normal Equations"

Review of Summation Rules

6.5

$a = \text{constant}$

$$* \sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$$

$$* \sum_{i=1}^n ax_i = ax_1 + ax_2 + \dots + ax_n = a \cdot \sum_{i=1}^n x_i$$

$$* \sum_{i=1}^n a = a + a + \dots + a = n \cdot a$$

$$* \sum_{i=1}^n x_i - \sum_{i=1}^n y_i = \sum_{i=1}^n (x_i - y_i)$$

Derivation of LS Estimates - Cont's

* "Normal Equations"

$$(1) \sum_{i=1}^n y_i = nb_0 + b_1 \sum x_i$$

$$(2) \sum_{i=1}^n x_i y_i = b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2$$

$$* \text{ Let } \bar{y} = \frac{1}{n} \sum y_i \quad \Rightarrow \quad \sum_{i=1}^n y_i = n\bar{y}$$

$$\bar{x} = \frac{1}{n} \sum x_i \quad \Rightarrow \quad \sum_{i=1}^n x_i = n\bar{x}$$

* Using this we have:

$$(1) b_0 = \bar{y} - b_1 \bar{x}$$

$$(2) \sum_{i=1}^n x_i y_i = b_0 \cdot n\bar{x} + b_1 \sum_{i=1}^n x_i^2$$

$$\Rightarrow \sum_{i=1}^n x_i y_i = (\bar{y} - b_1 \bar{x}) \cdot n\bar{x} + b_1 \sum_{i=1}^n x_i^2$$

$$b_1 = \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} = \frac{S_{xy}}{S_{xx}}$$

S_{xy} = sample covariance

$S_{xx} = (n-1)S_x^2$ it's $(n-1) \times$ sample variance of x_i 's

* Can show that:

$$b_1 = \frac{\sum (x_i - \bar{x}) \cdot (y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}} \quad (*)$$

* Exercises:

(1) Show that (*) is the same as the above results from normal equations.

(2) Show that b_0 and b_1 give minimum.

Properties of fitted line:

(1) Residuals:
$$e_i = y_i - \underbrace{(b_0 + b_1 x_i)}_{\hat{y}_i}$$
$$= y_i - (\bar{y} - b_1 \bar{x} + b_1 x_i)$$
$$= y_i - \bar{y} - b_1 (x_i - \bar{x})$$

$$\Rightarrow \sum_{i=1}^n e_i = \sum_{i=1}^n [y_i - \bar{y} - b_1 (x_i - \bar{x})]$$
$$= \sum_{i=1}^n (y_i - \bar{y}) - b_1 \sum_{i=1}^n (x_i - \bar{x}) = 0$$

$$\sum (y_i - \bar{y}) = 0$$

$$\sum (x_i - \bar{x}) = 0$$

why?

Properties of fitted line Cont's

(1) $\sum e_i = 0$

(2)
$$\begin{aligned} \sum_{i=1}^n \hat{y}_i &= \sum_{i=1}^n (b_0 + b_1 x_i) \\ &= \sum_{i=1}^n (\bar{y} - b_1 \bar{x} + b_1 x_i) \\ &= n\bar{y} - b_1 n\bar{x} + b_1 \underbrace{\sum_{i=1}^n x_i}_{n\bar{x}} \\ &= n\bar{y} - b_1 n\bar{x} + b_1 n\bar{x} \\ &= n\bar{y} = \sum_{i=1}^n y_i \end{aligned}$$

(3) $\sum_{i=1}^n e_i x_i = 0$

(4) $\sum_{i=1}^n e_i \hat{y}_i = 0$

Exercises!

Statistical Assumptions for SLR

(1) We assume that the linear model is appropriate, i.e., that indeed

$$E(Y|X) = \beta_0 + \beta_1 X$$

* $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad i = 1, \dots, n$

\downarrow r.v. \swarrow parameters \searrow fixed \searrow random

* So the additional assumptions are:

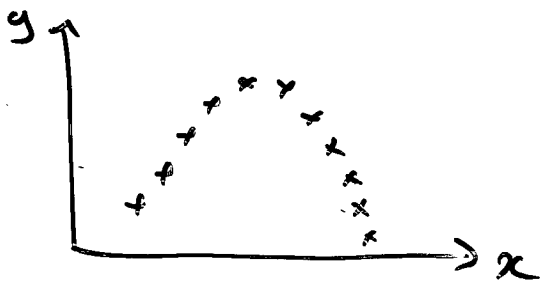
(2) $E(\epsilon_i) = 0 \implies E(Y_i | X_i = x_i) = \beta_0 + \beta_1 x_i$

(3) $Var(\epsilon_i) = \sigma^2 \implies Var(Y_i | X_i = x_i) = \sigma^2$

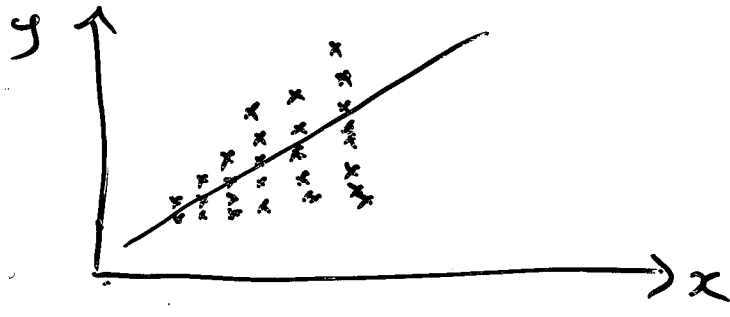
(4) ϵ_i 's are uncorrelated $\implies Y_i$'s are uncorrelated

Possible Violation of Statistical Assumptions

* linear model is not appropriate, for example



* Variance of error, $[var(y_i)]$ is increasing with x_i



* Can have influential points that moves the line.

