\* Announcement
① HW13 due 12/7
② Quiz 11 : Ch.11-2 ~ Ch12-2
③ RRR week teaching schedule
  • Tue : Topical OH
  • Wed : OH (2~4pm), Exam Walk-through (fall 2019)
  • Thu : GSI review sessions
        Mock Exam (Sp 2020)
  • Fri : Exam Walk-through
        (Sp 2020)

# STAT 88: Lecture 38

HW13 Part II Q3 (Ch11.4)
        Q4 (Ch12.1, 12.2)
          today and wed.

## Contents
Section 12.2: The Distribution of the Estimated Slope

Warm up:    Connecting Ch11 - Ch12.

Let $(X, Y)$ be a random pair and we observe $(x_1, Y_1), \ldots, (x_n, Y_n)$ from the linear regression model

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

want to estimate $E(X)$ from $X_1, \ldots, X_n$.
→ $\frac{1}{n}\sum_{i=1}^{n} X_i = \bar{X}$

(a) We defined the mean squared error of a linear function of $X$ as

$$\mathrm{MSE}(a, b) = E((Y - (aX + b))^2). \quad \text{or population quantity}$$

How would you estimate $\mathrm{MSE}(a, b)$ from $(x_1, Y_1), \ldots, (x_n, Y_n)$ ?

$$\frac{1}{n}\sum_{i=1}^{n}(Y_i - (ax_i + b)) = \widehat{\mathrm{MSE}}(a,b) \quad \xrightarrow{n \to \infty} \mathrm{MSE}(a,b)$$

(b) How can you estimate the best regression line, $\widehat{Y}_i = \widehat{\beta_0} + \widehat{\beta_1} x_i$?

Minimize $\widehat{\mathrm{MSE}}(a,b)$ w.r.t. $a$ and $b$

$\Rightarrow \quad a = \widehat{\beta_1} = \dfrac{\frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})(Y_i - \bar{Y})}{\frac{1}{n}\sum_{i=1}^{n}(x_n - \bar{x})^2}$ , $b = \widehat{\beta_0} = \bar{Y} - \widehat{\beta_1}\bar{x}$
(Exercise)

minimizes $\widehat{\mathrm{MSE}}(a,b)$

① First fix $a$ and minimize $\widehat{\mathrm{MSE}}(a,b)$ w.r.t. $b$
② Minimize $\widehat{\mathrm{MSE}}(a, b_a)$ w.r.t. $a$

(c) Compare the regression line in (b) with the population regression line $\widehat{Y} = \widehat{a}X + \widehat{b}$.

$\begin{cases} \widehat{a} = r\dfrac{\sigma_Y}{\sigma_X} = \dfrac{E((X-\mu_X)(Y-\mu_Y))}{\sigma_X^2} = \dfrac{E((X-\mu_X)(Y-\mu_Y))}{E((X-\mu_X)^2)} \\ \widehat{b} = \mu_Y - \widehat{a}\mu_X \end{cases}$

**Last time**

The simple regression model

$$Y = \text{response} \quad \text{and} \quad x = \text{predictor variable/covariate/feature}$$

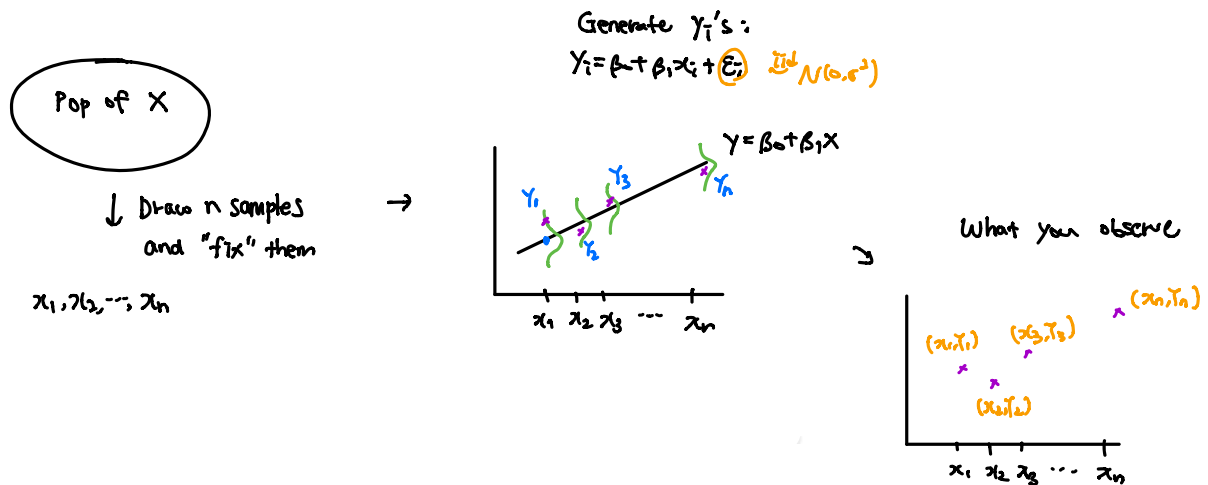<span style="color:blue">random</span>       <span style="color:blue">fixed</span>

We assume for each of $n$ observations

$$Y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\text{signal}} + \underbrace{\epsilon_i}_{\text{noise}},$$

$$\text{Var}(Y_i) = \text{Var}(\beta_0 + \beta_1 x_i + \epsilon_i)$$
$$= \text{Var}(\epsilon_i)$$
$$= \sigma^2$$

where

$$\sim \mathcal{N}(\beta_0 + \beta_1 x_i, \sigma^2)$$

- $\beta_0$ and $\beta_1$ are ==unobservable constant parameters==.

- $x_i$ is the value of the predictor variable for individual $i$ and ==is assumed to be constant== (that is, ==not random==).

- The errors $\epsilon_1, \epsilon_2, \ldots, \epsilon_n$ are ==i.i.d. normal $\mathcal{N}(0, \sigma^2)$ random variables==.

- The error variance $\sigma^2$ is an ==unobservable constant parameter==, and is assumed to be the ==same for all individuals== $i$.

Generate $Y_i$'s:
$$Y_i = \beta_0 + \beta_1 x_i + \underbrace{\epsilon_i}_{} \overset{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

Pop of X

$\downarrow$ Draw $n$ samples and "fix" them    $\rightarrow$

$x_1, x_2, \ldots, x_n$

$y = \beta_0 + \beta_1 x$

$Y_1$   $Y_3$   $Y_n$
   $Y_2$

$x_1 \ x_2 \ x_3 \ \cdots \ x_n$

$\searrow$

What you observe

$(x_n, Y_n)$
$(x_1, Y_1)$   $(x_3, Y_3)$
$(x_2, Y_2)$

$x_1 \ x_2 \ x_3 \ \cdots \ x_n$

Goal: Estimate $\beta_0$ and $\beta_1$ from $(x_1, Y_1), (x_2, Y_2), \ldots, (x_n, Y_n)$

$$\hat{\beta}_0, \hat{\beta}_1 \quad \rightarrow \quad \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

$\curvearrowleft$ random

## 12.2. The Distribution of the Estimated Slope

**Estimated Slope**

The least-squares estimate of the true slope $\beta_1$ is the slope of the regression line, given by

$$\widehat{\beta}_1 = \frac{\frac{1}{n}\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\frac{1}{n}\sum_{i=1}^n (x_i - \bar{x})^2}.$$ *random*

Is $\widehat{\beta}_1$ random or constant? What distribution does $\widehat{\beta}_1$ follow?
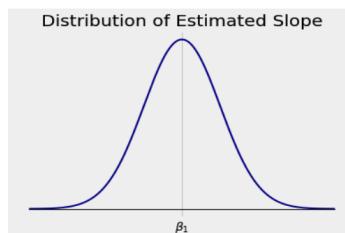
*Normal*

**Expectation of the Estimated Slope**

Let's find $E(\widehat{\beta}_1)$.

First, what is $E(Y_i)$ and $E(\bar{Y})$?

$\beta_0 + \beta_1 x_i$

$E(\bar{Y}) = \frac{1}{n}\sum_{i=1}^n E(Y_i) = \frac{1}{n}\sum_{i=1}^n (\beta_0 + \beta_1 x_i) = \beta_0 + \beta_1 \bar{x}$

$E(\widehat{\beta}_1) = \dfrac{\sum_{i=1}^n (x_i - \bar{x})\cdot E(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$

$\rightarrow \quad E(Y_i - \bar{Y}) = E(Y_i) - E(\bar{Y})$

$\qquad\qquad = \beta_1 (x_i - \bar{x})$

$\qquad = \dfrac{\sum_{i=1}^n (x_i - \bar{x})\cdot \beta_1 (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$

$\qquad = \beta_1$

Hence $\widehat{\beta}_1$ is an <mark>unbiased estimator</mark> of $\beta_1$.



Distribution of Estimated Slope

**Variance of the Estimated Slope**

FACT:
$$\mathrm{Var}(\widehat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}.$$

Note that as $n \to \infty$, $\sum_{i=1}^{n}(x_i - \bar{x})^2$ gets very large and $\mathrm{Var}(\widehat{\beta}_1) \to 0$ so the difference between $\widehat{\beta}_1$ and $\beta_1$ becomes very small with high probability. Hence
$$\widehat{\beta}_1 \sim \mathcal{N}\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right).$$

Example: (Exercise 12.4.1) Recall that the intercept of the regression line is given by the average of $Y$ minus the slope times the average of $x$. That is, $\widehat{\beta}_0 = \bar{Y} - \widehat{\beta}_1\bar{x}$. Is $\widehat{\beta}_0$ an unbiased estimate of $\beta_0$?

$\widehat{\beta}_0 = \bar{Y} - \widehat{\beta}_1 \bar{x}$

$E(\widehat{\beta}_0) = E(\bar{Y} - \widehat{\beta}_1 \bar{x})$

$\quad = E(\bar{Y}) - \bar{x} \, E(\widehat{\beta}_1)$

$\quad = \beta_0 + \beta_1 \bar{x} - \bar{x} \cdot \beta_1$

$\quad = \beta_0$

**Standard Error of the Estimated Slope**

$\varepsilon_i \sim N(0, \sigma^2) , \quad \sigma = SD(\varepsilon_i)$

$Y_i - (\beta_0 + \beta_1 x_i)$

$\uparrow$ Approximate.

$Y_i - (\widehat{\beta}_0 + \widehat{\beta}_1 x_i)$

$Y_i - \widehat{Y}_i = D_i$ "residual"

We have
$$\mathrm{SD}(\widehat{\beta}_1) = \frac{\sigma}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}}.$$

$\sigma$ is unknown so we have to estimate it. Recall $\sigma$ is the SD of the error, $\mathrm{SD}(\epsilon_1) = \sigma$. So we estimate $\sigma$ with the SD of the residuals. If
$$D_i = Y_i - \widehat{Y}_i = Y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i,$$

then
$$\widehat{\sigma} = \mathrm{SD}(D_1, \ldots, D_n) = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(D_i - \bar{D})^2}.$$

Sample standard deviation of $D_i$'s

This can be calculated in Python.

When the SD of an estimator is approximated by the data, it is called the SE (standard error):

$$\text{SE}(\widehat{\beta}_1) = \frac{\widehat{\sigma}}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}}.$$

When $n$ is large, it can be shown $\text{SE}(\widehat{\beta}_1)$ converges to $\text{SD}(\widehat{\beta}_1)$ so

$$T = \frac{\widehat{\beta}_1 - \beta_1}{\text{SE}(\widehat{\beta}_1)}$$

is approximately $\mathcal{N}(0, 1)$ for large $n$.

*To standardization of $\widehat{\beta}_1$*

$\rightarrow T_0 = \frac{\widehat{\beta} - \beta_1}{SD(\widehat{\beta}_1)} \sim \mathcal{N}(0,1)$

*Approximate $T_0$ by $T$*

## Pulse Rates

*Y*    *x*

We wish to predict active pulse rates from resting pulse rates.

pulse

*Y*    *x*

| Active | Rest | Smoke | Sex | Exercise | Hgt | Wgt |
|--------|------|-------|-----|----------|-----|-----|
| 97 | 78 | 0 | 1 | 1 | 63 | 119 |
| 82 | 68 | 1 | 0 | 3 | 70 | 225 |
| 88 | 62 | 0 | 0 | 3 | 72 | 175 |
| 106 | 74 | 0 | 0 | 3 | 72 | 170 |
| 78 | 63 | 0 | 1 | 3 | 67 | 125 |
| 109 | 65 | 0 | 0 | 3 | 74 | 188 |
| 66 | 43 | 0 | 1 | 3 | 67 | 140 |
| 68 | 65 | 0 | 0 | 3 | 70 | 200 |
| 100 | 63 | 0 | 0 | 1 | 70 | 165 |
| 70 | 59 | 0 | 1 | 2 | 65 | 115 |

... (222 rows omitted)

pulse.scatter('Rest', 'Active')



*Assumptions of simple linear regression model?*

```
active = pulse.column(0)
resting = pulse.column(1)
```

```
stats.linregress(x=resting, y=active)
```

Output:

$\widehat{\beta_1}$    (1.142879681904831,
$\widehat{\beta_0}$    13.182572776013345,
r    0.6041870881060092,
p-val   1.7861044071652305e-24,
SE($\widehat{\beta_1}$)   0.09938884436389145)

$n = 232$ is large so

Approximately

$$T = \frac{\widehat{\beta}_1 - \beta_1}{\mathrm{SE}(\widehat{\beta}_1)} \sim \mathcal{N}(0,1).$$

SD($\widehat{\beta_1}$)

A 95% CI for $\beta_1$ is

$$(\widehat{\beta}_1 \pm 2 \cdot \mathrm{SE}(\widehat{\beta}_1)) = (0.944, 1.342).$$

A fundamentally important question is whether the true slope $\beta_1$ is 0. If it is 0, then the resting pulse rate isn't involved in the prediction of the active pulse rate, according to the regression model. Our testing problem is

$$H_0 : \beta_1 = 0 \text{ vs } H_A : \beta_1 \neq 0.$$

$T$ is our test statistic. Under $H_0$,

$$T = \frac{\widehat{\beta}_1}{\mathrm{SE}(\widehat{\beta}_1)} \sim \mathcal{N}(0,1).$$

The observed value of the test statistic is 11.5. So the p-value is

$$\text{p-value} = P(T \geq 11.5) + P(T \leq -11.5) \approx 0.$$

We reject $H_0$ at 5% level.