# Stat 231

Nov 18, 2016

a

Tutorial Quiz: 30[th]

Syllabus — end of today

---

# Roadmap

* 5 min recap of SLRM.

* Least Square Estimates of and their properties

* Confidence Intervals / Testing of Hypotheses for $\alpha, \beta, \sigma$.

# Set-up:

$Y$ = response variable (trying to "explain" the variability of $Y$).

$X$ = explanatory variable. (given, not random)

& Objective: To try to estimate the relationship between $X$ and $Y$.

Based on ?

$(x_1, y_1)), \ldots \quad (x_n, y_n) \longrightarrow$ Sample we have drawn.

We make assumptions on how $X$ and $Y$ are related.

## SLRM:

$$Y_i \sim G_i\left(\alpha + \beta X, \sigma^2\right)$$

$$i = 1, \ldots n.$$

The mean response is a linear function of the explanatory variable.

Standard Gaussian model $\quad \mu(x) = \alpha + \beta x$

$$Y_i \sim G_i(\mu, \sigma)$$

$$i = 1, \ldots_n$$

# Notes:

$$Y_i \sim G(\mu, \sigma)$$

$$\Updownarrow$$

$$\boxed{Y_i = \mu + R_i,} \quad R_i \sim G(0, \sigma)$$

DETERMINISTIC PART

RANDOM PART

degrees of freedom = $n - \#$ of unknowns in the determinish part.
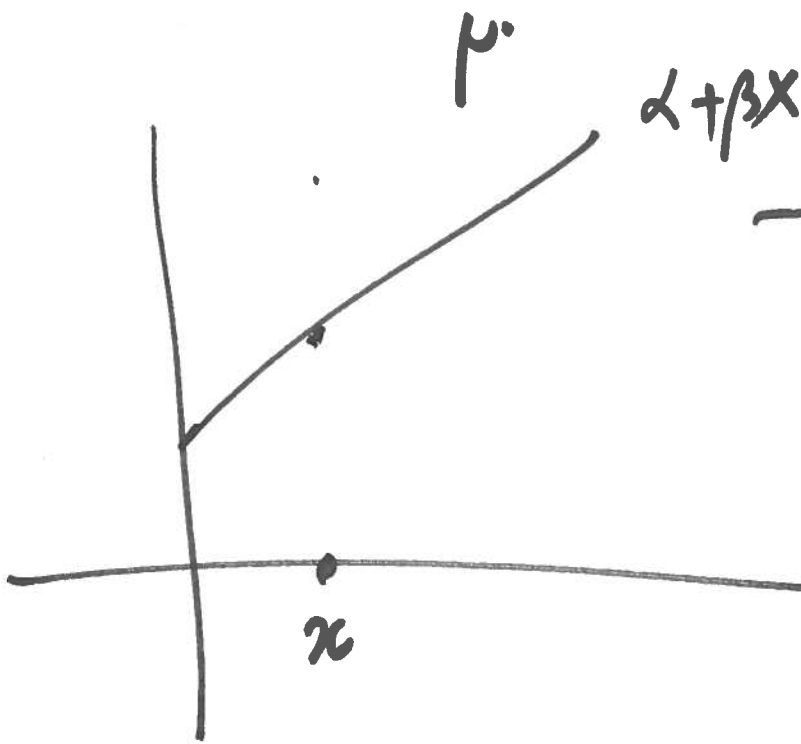
# Regression model

$$Y_i \sim G(\alpha + \beta x_i, \; \sigma)$$

$$\underbrace{\alpha + \beta x_i}_{\mu}$$

$$i = 1, \ldots n.$$

$$\Downarrow$$

$$Y_i = \underbrace{\alpha + \beta x_i}_{\mu} + R_i$$

$$R_i \sim G(0, \sigma)$$



$$\mu = \alpha + \beta x$$

Relevant
$$df = n - 2$$

# ESTIMATES

Method of Least Squares          Method of
                                 Max. Likelihood.

The estimates of $\alpha$ and $\beta$ are the same
for both methods.

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

$$\hat{\beta} = S_{xy}/S_{xx}$$

Least Square Estimates

$$\hat{\sigma}^2 = \frac{1}{n}\left[S_{yy} - \hat{\beta}S_{xy}\right]$$

$$S_{xx} = \sum (x_i - \bar{x})^2$$

$$S_{yy} = \sum (y_i - \bar{y})^2$$

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

---

Given a data set, we can always estimate $\hat{\alpha}, \hat{\beta}, \hat{\sigma}$ ($\hat{\beta}$ should be estimated first)

$$\hat{\sigma}^2 = \boxed{S^2 = \frac{1}{n-2} \left[ S_{yy} - \hat{\beta} S_{xy} \right]}$$

Sample Variance:

The Least Square Equation:

$$\hat{y} = \hat{\alpha} + \hat{\beta} x.$$

(best estimate for the linear relationship between X and Y).

$$\hat{\alpha} = \hat{\beta} \qquad \hat{\beta} = 0.9 \atop \hat{\alpha} = 10 \left.\vphantom{\begin{matrix}a\\a\\a\end{matrix}}\right\}$$

$$y = 10 + 0.9 X$$

$y = $ STAT 231

$x = $ STAT 230

# Confidence Interval for $\beta$.

## Fact 1 :

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{S_{xx}}$$

$$. = \frac{\sum(x_i - \bar{x})y_i}{S_{xx}} = \frac{\sum x_i(y_i - \bar{y})}{S_{xx}}$$

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

---

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

$$- \sum(x_i - \bar{x})\bar{y}$$

$$- \bar{y}\left(\sum(x_i - \bar{x})\right)$$

$$\sum (x_i - \bar{x})(y_i - \bar{y})$$

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

$$= \sum (x_i - \bar{x}) y_i - \bar{y} \left[ \sum x_i - \sum \bar{x} \right]$$

$$- \bar{y} \left[ n\bar{x} - n\bar{x} \right]$$

$$0$$

$$\hat{\beta} = \sum a_i y_i.$$

where $a_i = \dfrac{(x_i - \bar{x})^v}{S_{xy}}$

$$\boxed{\hat{\beta} = \sum_{i=1}^{n} a_i y_i.}$$

$a_i = \dfrac{x_i - \bar{x}}{S_{xx}}$

Each individual $\hat{\beta}$ can be thought of as outcome of a r.v. $\tilde{\beta}$.

What is the distribution of $\hat{\beta}$?

$$\tilde{\beta} = \sum a_i \cdot Y_i$$

Since $Y_i$'s are all Gaussian,

$\tilde{\beta}$s are also Gaussian.

---

**Theorem:**

$$\boxed{\tilde{\beta} \sim G\left(\beta, \frac{\sigma}{\sqrt{S_{xx}}}\right)}$$

$$\frac{\tilde{\beta} - \beta}{\dfrac{\sigma}{\sqrt{S_{xx}}}} \doteq Z$$

$$\boxed{\dfrac{\tilde{\beta} - \beta}{\dfrac{S}{\sqrt{S_{xx}}}} \sim T_{n-2}}$$

$n = 20$

$95\%$

(Row = 18
Column = 0.975)

Step 1 : To find $t^*$

Step 2 $P\left(-t^* < T < t^*\right) = 0.95$

$P\left(-t^* < \dfrac{\tilde{\beta} - \beta}{S/\sqrt{S_{xx}}} < t^*\right) = 0.95$

Coverage Interval:

$$\left( \tilde{\beta} \pm t^{\alpha} \frac{s}{\sqrt{Sxx}} \right)$$

Confidence Interval

$$\boxed{ \hat{\beta} \pm t^{\alpha} \frac{s}{\sqrt{Sxx}} }$$

Confidence Interval for $\beta$

$$H_0: \beta = \beta_0$$

$$D = \left| \frac{\hat{\beta} - \beta_0}{S/\sqrt{S_{xx}}} \right| \qquad \text{Test Statistic}$$

Calculate
$$d = \left| \frac{\hat{\beta} - \beta_0}{\frac{s}{\sqrt{S_{xx}}}} \right|$$

p-value:
$$P(D \geqslant d)$$
$$= P\left( |T_{n-2}| \geqslant d \right)$$

$$(x_1, y_1), \ldots \quad (x_n, y_n).$$

## Clicker question 1

All regression lines. must pass through $(\bar{x}, \bar{y})$

(a) True          61%

(b) False.          39%

$$y = \hat{\alpha} + \hat{\beta} x.$$

$$= \hat{\alpha} + \hat{\beta} \bar{x}$$

$$= \bar{y} - \hat{\beta}\bar{x} + \hat{\beta}\bar{x}$$

$$\bar{y} = \hat{\alpha} + \hat{\beta} \bar{x}.$$

# Clicker Question 2

If $r_{xy} = 0$, then $\hat{\beta} = 0$

(a) True $\longrightarrow$ 74%.

(b) False.