Test 2 Info. on Nexus

In Ch. 13, § 13.1

Last lecture:

Data:

$(x)$: Age (in years)

$(y)$: Price (in hundreds of dollars)

See Summary Stats. in Lecture 17.

a) least squares regression line:

$$\widehat{Price} = 318.83 - 33.49 \times Age$$

b) Give a brief interpretation of "a" and "b" Calculated in part (a).

b: $-33.49$ is the estimated decrease ($b < 0$)

in the mean price when age(x) is increased one unit (1 year).

That is, the average price of the vehicle goes down by $3349 when the age of the vehicle goes up by one year.

a: Statistically, 318.83 is the mean price when Age(x) is Zero. For this problem, 318.83 is not meaningful.

(C) Predict the price of a 7-year old car.

$$\widehat{Price} = 318.83 - 33.49 \times Age$$

$$\widehat{Price} = 318.83 - 33.49 \times 7 = 84.4$$

( i.e.   84.4 × 100  =  $ 8440 ).
                ↑
            units of x

Remark: It is safe to use the estimated

regression line to predict the price of a 7-year old car because $x=7$ falls within the range of $x$-values observed in this data set (See Lecture 17).

d) Predict the Price of an 18-year old car. Comment on the finding.

$$\widehat{Price} = 318.83 - 33.49 \times 18 = -283.99$$

i) negative price!

ii) Even if the Predicted Price for an 18-year old car were positive, $x=18$ is outside the range of $x$-values for which the data is collected; it is not safe to use the regression line for $x=18$.

Regression is for interpolation, not extrapolation.

§13.2 Standard Deviation of Errors & Coefficient of Determination:

Recall the idea of a residual:

$$e = Y - \hat{Y} \qquad \text{where,}$$

$Y$ = Observed response,

$\hat{Y}$ = Predicted response from the estimated least squares regression line.

Recall: The Simple Linear Regression Model:

$$Y_i = A + BX_i + \mathcal{E}_i$$

$$\mathcal{E}_i \overset{ind}{\sim} N(0, \sigma^2), \quad \text{for} \quad i = 1, \ldots, n.$$

$e_i$ is an estimate of $\mathcal{E}_i$.

Long-hand formula for estimating $\sigma$:

$$S_e = \hat{\sigma} = \sqrt{\frac{SSE}{n-2}}$$

$$= \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} = \sqrt{\frac{\sum e_i^2}{n-2}}$$

$$Se = \sqrt{\frac{SS_{yy} - b\, S_{xy}}{n-2}}$$

where

$$SS_{xy} = \sum xy - \frac{(\sum x)(\sum y)}{n}$$

$$SS_{yy} = \sum y^2 - \frac{(\sum y)^2}{n}$$

Ex. Consider the food expenditure (y) and monthly income (x) example in a previous lecture.

Summary Statistics:

$\sum x = 212, \quad \sum y = 64, \quad \sum xy = 2150,$

$\sum x^2 = 7,222, \quad \sum y^2 = \ldots \ldots ,$

$$SS_{yy} = \sum y^2 - \frac{(\sum y)^2}{n} = 60.8571$$

Q'n: What is an estimate of the Standard deviation of the residuals?

( Model: Expenditure = $A + B*Income + \varepsilon$

where $\varepsilon \overset{ind}{\sim} N(0, \sigma^2)$.

This q'n is asking you to estimate $\sigma$.

$$S_e = \sqrt{\frac{SS_{yy} - bS_{xx}}{n-2}}$$

$$= \sqrt{\frac{60.8571 - 0.2642(211.7143)}{7-2}} = 0.9922$$

$\therefore \widehat{\sigma}_e = 0.9922.$

[ Extra: STAT-3103

$$SST = SSR + SSE$$

$$SST = \sum_{i=1}^{n} (y_i - \bar{y})^2 \qquad Total\ Sums\ of\ Squares$$

$$SSR = \sum_{i=1}^{n} (\widehat{y_i} - \bar{y})^2 \qquad Regression\ Sums\ of\ Squares$$

$$SSE = \sum_{i=1}^{n} (y_i - \widehat{y_i})^2 \qquad Error\ Sums\ of\ Squares.$$

For the least squares line to be useful, we want SSR to be "large" relative to SSE. ]

A statistic to assess the usefulness of the regression line is the Coefficient of determination $(r^2)$.

It measures the proportion of the variation in Y (the response variable) that is explained by our X (the explanatory variable).

$$r^2 = \frac{SSR}{SST}$$

Short-cut formula:

$$r^2 = b \frac{SSxy}{SSyy}$$

( Not given ! )

But, to find r, square the value of

$$r = \frac{SS_{xy}}{\sqrt{SS_{xx}SS_{yy}}}$$

← Correlation between $x$ and $y$

↑ Given on Final

Property of $r^2$: $\qquad 0 \le r^2 \le 1$

Ex. Consider the monthly food expenditure ($y$) and monthly income ($x$) regression example.

Q'n: What is the percentage of the variation in monthly food expenditure that is explained by monthly income ($x$)?

(ie What is $r^2$ (expressed as a percentage)?)

$$r = \frac{SS_{xy}}{\sqrt{SS_{xx} \cdot SS_{yy}}} = \frac{211.7143}{\sqrt{801.4286\,(60.8571)}} = 0.9587 \approx 0.96$$

$$\therefore \quad r^2 = 0.96^2 = 0.92$$

*Interpretation:* Approx. 92% of the variation in monthly food expenditure (y) is explained by monthly income (x).

Compare with:

$$r^2 = \frac{b \, SSxy}{SSyy} = 0.2642 \times \frac{211.7143}{60.8571} = 0.92$$
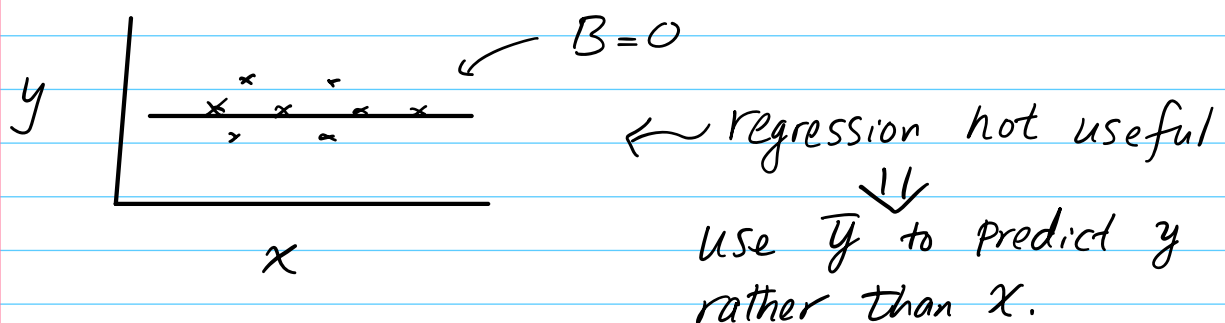
## § 13.3 Inference about B

*Idea:* Want to test whether the regression line is "useful". We do this by testing.

$H_0: B = 0$    "regression line not useful".

against

$H_1: B > 0$   or   $H_1: B < 0$   or

$H_1: B \neq 0$.



$B = 0$

← regression not useful

$\Downarrow$

use $\bar{y}$ to predict y rather than x.

Test Statistic:

$$t = \frac{b - B_0}{S_e / \sqrt{SS_{xx}}}$$

(see $S_e$ on formula sheet)

Under $H_0$ (ie. assuming $H_0$ is true),

$$t = \frac{b - B_0}{S_e / \sqrt{SS_{xx}}} \sim t_{n-2}$$

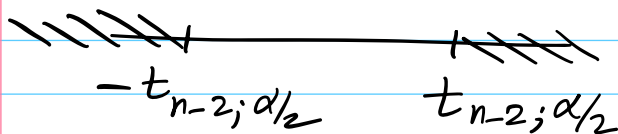$H_0 : B = 0$

| Test | | Rejection Region : |
|------|--|--------------------|

$H_0 : B = 0$ vs. $H_1 : B > 0$    $t > t_{n-2; \alpha}$

$H_0 : B = 0$ vs. $H_1 : B < 0$    $t < -t_{n-2; \alpha}$

$H_0 : B = 0$ vs. $H_1 : B \neq 0$    $|t| > t_{n-2; \alpha/2}$

$$\Leftrightarrow$$

$$t < -t_{n-2; \alpha/2} \quad \text{or} \quad t > t_{n-2; \alpha/2}$$

$-t_{n-2; \alpha/2} \qquad t_{n-2; \alpha/2}$

Omit p-value for testing $H_0: B = 0$.