

11 Lecture 11: Feb 13

Last time

- Introduction of simple linear regression

Today

- HW2 posted
- The statistical model of the SLR (JF chapter 6)
- Properties of the Least-Squares estimator
- Inference of SLR model

R-square

The ratio of RegSS to TSS is called the *coefficient of determination*, or sometimes, simply “r-square”. it represents the proportion of variation observed in the response variable y which can be “explained” by its linear association with x .

- In simple linear regression, “r-square” is in fact equal to r_{xy}^2 . (But this isn’t the case in multiple regression.)
- It is also equal to the squared correlation between y_i and \hat{y}_i . (This is the case in multiple regression.)

For Davis’s regression of measured on reported weight:

$$\text{TSS} = 4753.8$$

$$\text{RSS} = 418.87$$

$$\text{RegSS} = 4334.9$$

Thus,

$$r^2 = \frac{4334.9}{4753.8} = 1 - \frac{418.87}{4753.8} = 0.9119$$

The statistical model of Simple Linear Regression

Standard statistical inference in simple regression is based on a *statistical model* that describes the population or process that is sampled:

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

where the coefficients β_0 and β_1 are the *population regression parameters*. The data are randomly sampled from some population of interest.

- y_i is the value of the response variable

- x_i is the explanatory variable
- ϵ_i represents the aggregated omitted causes of y (i.e., the causes of y beyond the explanatory variable), other explanatory variables that could have been included in the regression model, measurement error in y , and whatever component of y is inherently random.

Key assumptions of SLR

The key assumptions of the SLR model concern the behavior of the errors, equivalently, the distribution of y conditional on x :

- *Linearity.* The expectation of the error given the value of x is 0: $\mathbf{E}(\epsilon) \equiv \mathbf{E}(\epsilon|x_i) = 0$. And equivalently, the expected value of the response variable is a linear function of the explanatory variable: $\mu_i \equiv \mathbf{E}(y_i) \equiv \mathbf{E}(y_i|x_i) = \mathbf{E}(\beta_0 + \beta_1 x_i + \epsilon_i|x_i) = \beta_0 + \beta_1 x_i$.
- *Constant variance.* The variance of the errors is the same regardless of the value of x : $\mathbf{Var}(\epsilon|x_i) = \sigma_\epsilon^2$. The constant error variance implies constant conditional variance of y on given x : $\mathbf{Var}(y|x_i) = \mathbf{E}((y_i - \mu_i)^2) = \mathbf{E}((y_i - \beta_0 - \beta_1 x_i)^2) = \mathbf{E}(\epsilon_i^2) = \sigma_\epsilon^2$. (Question: why the last equal sign?)
- *Normality.* The errors are independent identically distributed with Normal distribution with mean 0 and variance σ_ϵ^2 . Write as $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$. Equivalently, the conditional distribution of the response variable is normal: $y_i \stackrel{iid}{\sim} N(\beta_0 + \beta_1 x_i, \sigma_\epsilon^2)$.
- *Independence.* The observations are sampled independently.
- *Fixed X , or X measured without error and independent of the error.*
 - For experimental research where X values are under direct control of the researcher (i.e. X 's are fixed). If the experiment were replicated, then the values of X would remain the same.
 - For research where X values are sampled, we assume the explanatory variable is measured without error and the explanatory variable and the error are independent in the population from which the sample is drawn.
- *X is not invariant.* X 's can not be all the same.

Figure 11.1 shows the assumptions of linearity, constant variance, and normality in SLR model.

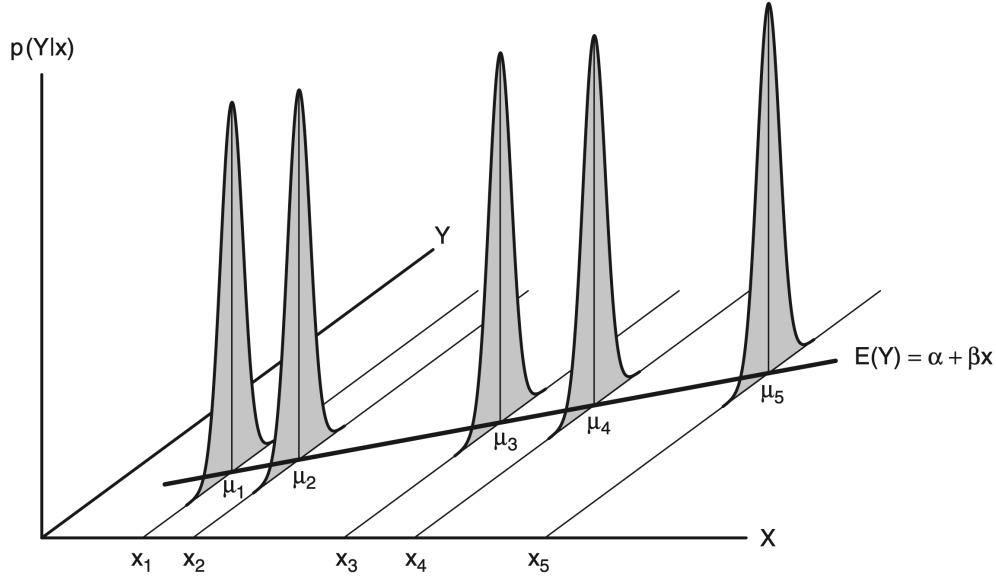


Figure 11.1: The assumptions of linearity, constant variance, and normality in simple regression. The graph shows the conditional population distributions $\Pr(Y|x)$ of Y for several values of the explanatory variable X , labeled as x_1, x_2, \dots, x_5 . The conditional means of Y given x are denoted μ_1, \dots, μ_5 .

Properties of the Least-Squares estimator

Under the strong assumptions of the simple regression model, the sample least squares coefficients $\hat{\beta}_{ls}$ have several desirable properties as estimators of the population regression coefficients β_0 and β_1 :

- The least-squares intercept and slope are *linear estimators*, in the sense that they are linear functions of the observations y_i .

Proof:

- The sample least-squares coefficients are *unbiased estimators* of the population regression coefficients:

$$\mathbf{E}(\hat{\beta}_0) = \beta_0$$

$$\mathbf{E}(\hat{\beta}_1) = \beta_1$$

Proof:

- Both $\hat{\beta}_0$ and $\hat{\beta}_1$ have simple sampling variances:

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}$$

Proof:

- Rewrite the formula for $\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{(n-1)S_X^2}$, we see that the sampling variance of the slope estimate will be small when
 - The error variance σ_ϵ^2 is small
 - The sample size n is large
 - The explanatory-variable values are spread out (i.e. have a large variance, S_X^2)
 - (Gauss-Markov theorem) Under the assumptions of linearity, constant variance, and independence, the least-squares estimators are BLUE (Best Linear Unbiased Estimator), that is they have the smallest sampling variance and are unbiased. (show this)
- Proof:*

- Under the full suite of assumptions, the least-squares coefficients $\hat{\beta}_0$ and $\hat{\beta}_1$ are the maximum-likelihood estimators of β_0 and β_1 . (show this)
- Proof:*

- Under the assumption of normality, the least-squares coefficients are themselves normally distributed. Summing up,

$$\begin{aligned}\hat{\beta}_0 &\sim N\left(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}\right) \\ \hat{\beta}_1 &\sim N\left(\beta_1, \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}\right)\end{aligned}$$

Statistical inference of the SLR model

Now we have the distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

$$\begin{aligned}\hat{\beta}_0 &\sim N\left(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}\right) \\ \hat{\beta}_1 &\sim N\left(\beta_1, \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}\right).\end{aligned}$$

However, σ_ϵ is never known in practice. Instead, an *unbiased* estimator of σ_ϵ^2 is given by

$$\hat{\sigma}_\epsilon^2 = MS[E] = \frac{SS[E]}{n-2}.$$

Proof:

Confidence intervals

Now we substitute $\hat{\sigma}_\epsilon^2$ into the distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

$$\begin{aligned}\hat{\beta}_1 &\sim N(\beta_1, \frac{\sigma_\epsilon^2}{\sum(x_i - \bar{x})^2}) \\ \hat{\beta}_0 &\sim N(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum(x_i - \bar{x})^2})\end{aligned}$$

to get the estimated standard errors:

$$\begin{aligned}\widehat{SE}(\hat{\beta}_1) &= \sqrt{\frac{MS[E]}{\sum(x_i - \bar{x})^2}} \\ \widehat{SE}(\hat{\beta}_0) &= \sqrt{MS[E] \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right)}\end{aligned}$$

And the $100(1 - \alpha)\%$ confidence intervals for β_1 and β_0 are given by

$$\begin{aligned}\hat{\beta}_1 \pm t(n - 2, \alpha/2) \sqrt{\frac{MS[E]}{S_{xx}}} \\ \hat{\beta}_0 \pm t(n - 2, \alpha/2) \sqrt{MS[E] \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}\end{aligned}$$

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

Confidence interval for $\mathbf{E}(Y|X = x_0)$

The conditional mean $\mathbf{E}(Y|X = x_0)$ can be estimated by evaluating the regression function $\mu(x_0)$ at the estimates $\hat{\beta}_0, \hat{\beta}_1$. The conditional variance of the expression isn't too difficult (already shown):

$$\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_0 | X = x_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

This leads to a confidence interval of the form

$$\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t(n - 2, \alpha/2) \sqrt{MS[E] \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}$$

Prediction interval

Often, prediction of the response variable Y for a given value, say x_0 , of the independent variable of interest. In order to make statements about future values of Y , we need to take into account

- the sampling distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

- the randomness of a future value Y .

We have seen the predicted value of Y based on the linear regression is given by $\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$.

The 95% prediction interval has the form

$$\hat{Y}_0 \pm t(n-2, \alpha/2) \sqrt{MS[E] \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}.$$

Hypothesis test

To test the hypothesis $\boxed{H_0 : \beta_1 = \beta_{slope_0}}$ that the population slope is equal to a specific value β_{slope_0} (most commonly, the null hypothesis has $\beta_{slope_0} = 0$), we calculate the test statistic (T -statistics) with $df = n - 2$

$$t_0 = \frac{\hat{\beta}_1 - \beta_{slope_0}}{\widehat{SE}(\hat{\beta}_1)} \sim t_{n-2}$$