

Supplementary Handout for Dis 13: Simple Linear Regression

1 Motivation

- Congratulation! You've made it to the last week of our Econ 310 discussion section.
- Throughout this semester, we have learned many useful statistical tools. But the most important tool is **statistical inference** – using statistics and their relevant distributions to draw inference about population parameters.
- Statistical inference is at the heart of your future econometrics class (either Econ 400 or 410).
- You might ask: what are we going to study in econometrics?
 - “Econometrics is the application of **statistical methods** to economic data in order to give empirical content to **economic relationships**.” – Wikipedia
(Put in other words, econometrics asks you to recover true economic relationships from sample data by using statistical methods.)
 - Some statistical methods used to describe relationship between variables:
 - * **Covariance**
→ Positive vs. negative relationship, but the scale of relationship is ambiguous.
 - * **Coefficient of correlation**
→ Positive vs. negative relationship, and the scale is normalized between -1 and 1 (inclusive). But unclear on how change in one variable quantifies to change in the other.
 - * **Simple linear regression**
→ Positive vs. negative relationship. Measurable scale (through statistical inference on slope coefficient). Tells us rate of change from x to y (via the slope coefficient β_1).

2 Simple Linear Regression

- Suppose that high school GPA (x) and student SAT test score (y) follow a linear relationship, which is expressed as the following:

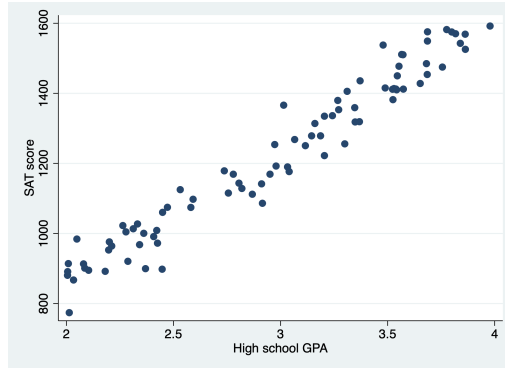
$$y = \beta_0 + \beta_1 x + \varepsilon$$

Notice that the above describes the **true** relationship between y and x . Some terminologies:

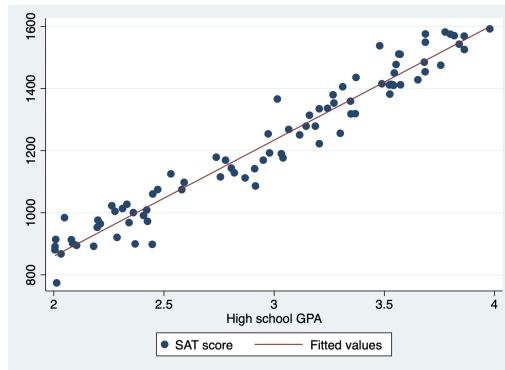
- y : dependent variable (in this case, the student SAT score)
- x : independent variable (in this case, the student high school GPA)
- β_0 : true intercept
- β_1 : true slope
- ε : error term
- With a true linear relationship established, how do we recover the true parameters β_0 and β_1 ?
⇒ use sample data to estimate them!

2.1 How to estimate the linear line?

- Let's say that we collect a simple random sample of size n : $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$. Again, x is high school GPA, and y is the student's corresponding SAT test score.
- The scatter plot created for the sample data looks like the following:



- How to fit a linear line in our sample data in order to estimate β_0 and β_1 ?



Solution: minimize the **sum of squared distance** between each (x_i, y_i) data point and the fitted line.

- Denote the estimated fitted line as \hat{y} , where

$$\hat{y} = b_0 + b_1x$$

- \hat{y} : predicted y / fitted y
- b_0 : estimate of intercept
- b_1 : estimate of slope

The estimates of b_0 and b_1 are obtained by

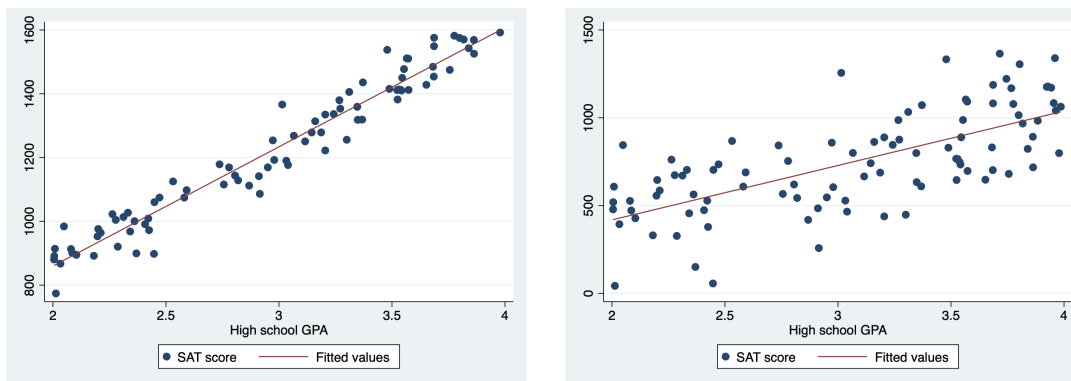
$$\min_{b_0, b_1} \sum_i (y_i - \hat{y}_i)^2 \quad \Leftrightarrow \quad \min_{b_0, b_1} \sum_i (y_i - b_0 - b_1x_i)^2$$

Taking first order conditions to solve the minimization problem yields

$$b_1 = \frac{s_{xy}}{s_x^2} \quad b_0 = \bar{y} - b_1\bar{x}$$

2.2 How to evaluate level of fitness for the estimated linear line?

- In the scatter plot that we just looked at, the estimated linear line seems to fit the data pretty well.
- However, say that you've gathered another set of data on high school GPA (x) and SAT score (y), but this new dataset is a bit more noisy. The scatter plot of the new dataset looks like the one on the right:



- How can we tell the above two cases apart?
 - The estimated line fits the data on the left better than the data on the right.
 - Use some goodness of fit measure to differentiate the level of fitness for the estimated line!
- Goodness of fit measure: **coefficient of determination**, or R^2

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

where

- SSE = sum of squares for error = $\sum_{i=1}^n (y_i - \hat{y}_i)^2$
- SST = total sum of squares = $\sum_{i=1}^n (y_i - \bar{y})^2$
- Some properties of R^2 :
 - $0 \leq R^2 \leq 1$
 - The better fit the projected line, the bigger the R^2 .
 - **Interpretation of R^2 :**
Say that $R^2 = 0.35$, then about 35% of the sample variation in y is explained by variation in x .

2.3 How to interpret the linear line?

- Continue with the example of student high school GPA (x) and SAT score (y). Say that the fitted line is estimated to be $\hat{y} = 120 + 400x$
- How do we interpret $b_1 = 400$?
 - b_1 gives us an estimate on the rate of change from x onto y .
 - **For this example:** For every 1.0 point increase in high school GPA, their predicted SAT score increases by 400.

- How do we interpret $b_0 = 120$?
 - b_0 tells us what predicted y would be if $x = 0$.
 - **For this example:** If a student has high school GPA = 0, then their predicted SAT score is 120.
 - Keep in mind, if your dataset doesn't contain actual data points near where $x = 0$, then it might not make sense for you to interpret b_0 .
For example, it is virtually impossible for us to observe a student in our data where their GPA is 0.0. So saying that someone with 0.0 GPA has predicted SAT score equals to 120 (a) doesn't make sense, and (b) doesn't really offer any meaningful information.

2.4 How to perform hypothesis testing on the slope coefficient?

- Since the estimate of β_1 (i.e. b_1) gives us an estimate on the rate of change from x onto y , it is often informative for one to perform hypothesis testing on the true slope coefficient (often to check whether β_1 is positive / negative / equal to 0).
- If one goes down the test statistic & rejection region route:

$$\text{test statistic} = \frac{b_1 - \beta_{1,H_0}}{s_\varepsilon / \sqrt{(n-1)s_x^2}} \sim t_{n-2}$$

where s_ε is the standard error of the regression:

$$s_\varepsilon = \text{Root MSE (Root Mean Square Error)} = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}}$$

- If one goes down the confidence interval route, the confidence interval of $(1 - \alpha)$ confidence level is

$$\left[b_1 - t_{\alpha/2} \times \frac{s_\varepsilon}{\sqrt{(n-1)s_x^2}}, \quad b_1 + t_{\alpha/2} \times \frac{s_\varepsilon}{\sqrt{(n-1)s_x^2}} \right]$$

Side note on s_ε (Root MSE):

- s_ε measures the standard deviation of regression:
 - If the linear line (i.e. the linear model) fits the data well, then s_ε should be relatively small.
 - If the linear line fits the data poorly, then s_ε should be relatively large.
- s_ε can be found from Stata regression output:

. reg Y X						
Source	SS	df	MS	Number of obs	=	85
Model	4151960.23	1	4151960.23	F(1, 83)	=	1577.37
Residual	218472.749	83	2632.2018	Prob > F	=	0.0000
				R-squared	=	0.9500
				Adj R-squared	=	0.9494
Total	4370432.98	84	52028.9641	Root MSE	=	51.305

In the upper left corner, the *SS* column, *Residual* row records the value of SSE, which in this specific example is 218472.749.

The *Number of obs* = 85 on the right hand side tells us that $n = 85$. This allows us to calculate s_ϵ by hand:

$$s_\epsilon = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - 2}} = \sqrt{\frac{218472.749}{85 - 2}} = 51.305$$

Alternatively, looking at the right side, Root MSE = 51.305 directly tells us the value of s_ϵ .