

## Problem Set 2

Due by 3/7, 11:59PM

Consider a simple regression model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad (1)$$

for  $i = 1, 2, \dots, n$ . Stacking the observations, we can write the model for the entire sample as

$$\mathbf{y} = \beta_0 \mathbf{l}_n + \beta_1 \mathbf{x} + \boldsymbol{\varepsilon} \quad (2)$$

where  $\mathbf{l}_n$  is the  $n$  by 1 vector of 1's; and  $\mathbf{y}$ ,  $\mathbf{x}$ ,  $\boldsymbol{\varepsilon}$  are all  $n$  by 1 vectors. Here  $\{y_i, x_i\}$ 's are random observations (i.i.d.), and  $\beta_0$  and  $\beta_1$  are unknown coefficients. The least-squares methodology tries to infer about the unknown coefficients by finding a best fit to sample data taken from the population. The best fit refers to values of  $\beta_0$  and  $\beta_1$  so that the sum of the squared residuals (or prediction errors) is the least possible, i.e.

$$(\hat{\beta}_0, \hat{\beta}_1)' = \arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.$$

It can be shown that the least squares estimator for the slope coefficient  $\beta_1$  is given by

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

and the least-squares estimator for the intercept is given by

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

where  $\bar{y}$  and  $\bar{x}$  denote the sample means of  $y$  and  $x$ , respectively.

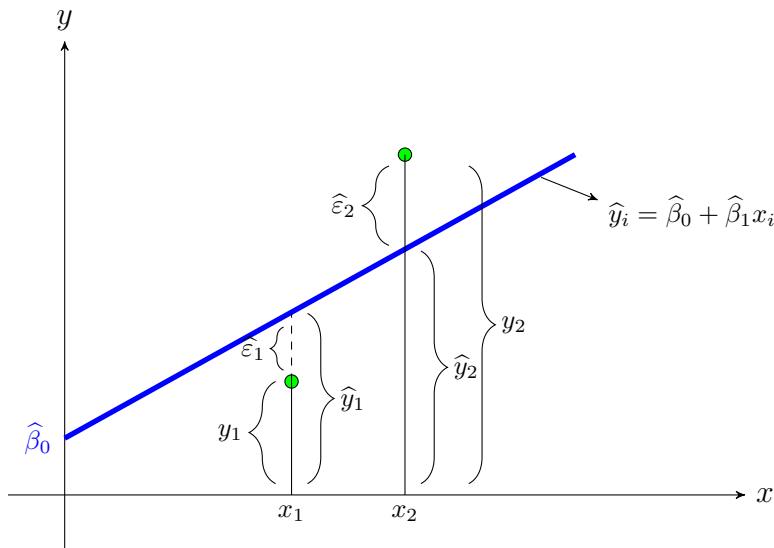


FIGURE 1. The least-squares regression line.

- (0) Write a function which will take  $\mathbf{y}$  and  $\mathbf{x}$  as inputs and return estimates of  $\beta_0$  and  $\beta_1$  using the least-squares methodology.
- (1) Set the seed to 37 for the random number generator in [numpy](#), i.e., `np.random.seed(37)`.

(2) Generate an array of zeros of size 5000, and name it `tmp1`. For 5000 times, do the following in a loop:

- (a) generate 1000 observations on  $\mathbf{x}$  by drawing randomly from the standard normal distribution, (i.e.,  $\mathbf{x}$  is a 1000 by 1 vector of standard Gaussian random variables);
- (b) generate 1000 observations on  $\boldsymbol{\varepsilon}$  by drawing randomly from the standard normal distribution, (i.e.,  $\boldsymbol{\varepsilon}$  is a 1000 by 1 vector of standard Gaussian random variables);
- (c) generate the  $\mathbf{y}$  vector using

$$\mathbf{y} = 0.5 \mathbf{l}_n + 1.8 \mathbf{x} + \boldsymbol{\varepsilon}$$

using the draws from (a) and (b);

- (d) estimate the simple linear regression model using your function;

- (e) save  $\beta_1$  estimate to `tmp1`.

(3) Calculate the mean of `tmp1`. Is it 1.8? Is it close to 1.8?

(4) Redo (2) and (3), but this time in 2(b), generate 1000 observations on  $\boldsymbol{\varepsilon}$  as

$$\varepsilon_i = -0.5x_i + v_i$$

for  $i = 1, 2, \dots, 1000$ , where  $v_i$  is drawn from the standard normal distribution.

(5) Calculate the mean of `tmp1` again. Compared to your finding in (3), is it further away from 1.8? Why?