

SIMPLE LINEAR REGRESSION

DAVAR KHOSHNEVISAN

1. LINE-FITTING

Points $(x_1, y_1), \dots, (x_n, y_n)$ are given; e.g. on a scatterplot.

What is “the best line” that describes the relationship between the x ’s and the y ’s?

To understand this better, let us focus on a line $L(x) = \alpha + \beta x$ where $\alpha, \beta \in \mathbf{R}$ are fixed but otherwise arbitrary.

“Fitting L to the points (x_i, y_i) ” means estimating y_i by $L(x_i)$ for all $i = 1, \dots, n$. The error, e_i , in estimating y_i by $L(x_i)$ is called the i^{th} residual.

We choose “the best line of fit” according to the *least squares principle* of C. Gauss:¹ Minimize $\sum_{i=1}^n e_i^2$ among all possible lines of the form $L(x) = \alpha + \beta x$. This is done by doing a little calculus: For a fixed line $L(x) = \alpha + \beta x$,

$$(1) \quad \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (L(x_i) - y_i)^2 = \sum_{i=1}^n (\alpha + \beta x_i - y_i)^2.$$

Call this $\mathcal{H}(\alpha, \beta)$. Then we are asked to find a and b that minimize \mathcal{H} . But this is a two-variable calculus problem. It turns out to be enough to solve:

$$(2) \quad \frac{\partial \mathcal{H}(\alpha, \beta)}{\partial \alpha} = 0 \quad \text{and} \quad \frac{\partial \mathcal{H}(\alpha, \beta)}{\partial \beta} = 0.$$

Date: August 30, 2004.

¹Sometimes, people have reason to use other, more “robust,” principles. A common alternative in such a case is the principle of “least absolute deviation.” It seeks to find a line that minimizes $\sum_{i=1}^n |e_i|$. Occasionally, this is also called the “ \mathcal{L}^1 method”; this is to distinguish it from least squares which is also sometimes called the “ \mathcal{L}^2 method.”

First,

$$\begin{aligned}
 \frac{\partial \mathcal{H}(\alpha, \beta)}{\partial \alpha} &= \sum_{i=1}^n 2(\alpha + \beta x_i - y_i) \\
 &= 2\alpha n + 2\beta \sum_{i=1}^n x_i - 2 \sum_{i=1}^n y_i, \text{ and} \\
 \frac{\partial \mathcal{H}(\alpha, \beta)}{\partial \beta} &= \sum_{i=1}^n 2(\alpha + \beta x_i - y_i)x_i \\
 &= 2\alpha \sum_{i=1}^n x_i + 2\beta \sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i y_i.
 \end{aligned}
 \tag{3}$$

These are called the *normal equations*. Now (2) is transformed into two equations in two unknowns $[\alpha$ and $\beta]$:

$$\begin{aligned}
 \alpha \bar{x} + \beta \overline{x^2} &= \overline{xy} \\
 \alpha + \beta \bar{x} &= \bar{y}.
 \end{aligned}
 \tag{4}$$

Multiply the second equation of (4) by \bar{x} , and then subtract from the first equation to find that $\beta(\overline{x^2} - (\bar{x})^2) = \overline{xy} - \bar{x} \cdot \bar{y}$. You should recognize this, in statistical terms, as $\beta \text{Var}(x) = \text{Cov}(x, y)$. Equivalently, $\beta = \text{Cov}(x, y) / \text{Var}(x) = \text{Corr}(x, y) \text{SD}_y / \text{SD}_x$. Plug this into the second equation of (4) to find that $\alpha = \bar{y} - \beta \bar{x}$ for the computed β . To summarize,

Theorem 1.1. *The least-squares line through $(x_1, y_1), \dots, (x_n, y_n)$ is unique and defined by*

$$L(x) = \bar{y} + \hat{\beta}(x - \bar{x}), \quad \text{where} \quad \hat{\beta} = \text{Corr}(x, y) \frac{\text{SD}_y}{\text{SD}_x}.
 \tag{5}$$

2. THE MEASUREMENT-ERROR MODEL

Let $\mathbf{Y} = (Y_1, \dots, Y_n)$ be a random sample of n i.i.d. copies of the response variable. The *measurement-error* model posits the following:

$$Y_i = \alpha + \beta X_i + \varepsilon_i \quad i = 1, \dots, n.
 \tag{6}$$

Here, $\mathbf{X} = (X_1, \dots, X_n)$ is a non-random vector of constants—the explanatory variables—and α and β are unknown parameters. This model also assumes that the ε_i 's are i.i.d. $N(0, \sigma^2)$ for an unknown parameter $\sigma > 0$. The Y_i 's are random only because the ε_i 's are (and not the X_i 's).

According to the principle of least squares (Theorem 1.1), the best least-squares estimates of α and β are, respectively,

$$\hat{\beta} = \frac{\frac{1}{n} \sum_{i=1}^n X_i Y_i - \bar{X} \cdot \bar{Y}}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} \quad \text{and} \quad \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}.
 \tag{7}$$

The more important parameter is β . For instance, consider the test,

$$H_0 : \beta = 0 \quad \text{vs.} \quad H_1 : \beta \neq 0.
 \tag{8}$$

This is testing the hypothesis that the explanatory variable X has a (linear) effect on Y .

So we need the distribution of $\hat{\beta}$. Note that

$$(9) \quad \hat{\beta} = \frac{\sum_{i=1}^n Y_i (X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \sum_{i=1}^n b_i Y_i,$$

where

$$(10) \quad b_i = \frac{X_i - \bar{X}}{ns_{\mathbf{X}}^2} \quad \text{for} \quad s_{\mathbf{X}}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Recall that the X_i 's are not random. Therefore, neither are the b_i 's. Also recall,

Lemma 2.1. *If V_1, \dots, V_n are independent and $V_i \sim N(\mu_i, \sigma_i^2)$, then for all non-random c_1, \dots, c_n , $\sum_{i=1}^n c_i V_i \sim N(\mu, \sigma^2)$ where $\mu = \mu_1 + \dots + \mu_n$ and $\sigma^2 = \sigma_1^2 + \dots + \sigma_n^2$.*

Consequently, $\hat{\beta} \sim N(\sum_{i=1}^n b_i E[Y_i], \sum_{i=1}^n b_i^2 \text{Var}(Y_i))$. But $Y_i = \alpha + \beta X_i + \varepsilon_i$. So, $E[Y_i] = \alpha + \beta X_i$, and $\text{Var}(Y_i) = \text{Var}(\varepsilon_i) = \sigma^2$. It is easy to check that: (i) $\sum_{i=1}^n b_i = 0$; (ii) $\sum_{i=1}^n b_i X_i = 1$; and (iii) $\sum_{i=1}^n b_i^2 = 1/(ns_{\mathbf{X}}^2)$. This proves that

$$(11) \quad \hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{ns_{\mathbf{X}}^2}\right).$$

Therefore, $E[\hat{\beta}] = \beta$. That is, $\hat{\beta}$ is an unbiased estimator of β . Moreover, if we knew σ^2 , then we could perform the test of hypothesis (8) at any prescribed level, say at 95%. The trouble is that we generally do not know σ^2 .

Because the Y_i 's have variance σ^2 , we can estimate σ^2 by $s_{\mathbf{Y}}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$. But then we need the joint distribution of $(\hat{\beta}, s_{\mathbf{Y}}^2)$. The key to this theory is that $\hat{\beta}$ is independent of $s_{\mathbf{Y}}^2$. We just determined the distribution of $\hat{\beta}$, and we will see later on that the H_0 -distribution of $s_{\mathbf{Y}}^2$ is essentially χ^2 . The rest will be smooth sailing.

To recap, we need to accomplish two things:

- (1) Derive the independence of $\hat{\beta}$ and $s_{\mathbf{Y}}^2$; and
- (2) Honestly compute the distribution of $s_{\mathbf{Y}}^2$ under H_0 .

Just about all of this semester's work is concerned with accomplishing these two goals (for more general models).