

SMF SOLUTIONS 7. 1.6.2012

Q1. To fit a straight line

$$y = a + bx$$

by least squares through a data set $(x_1, y_1), \dots, (x_n, y_n)$, we choose a, b so as to minimise

$$SS := \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a - bx_i)^2.$$

Taking $\partial SS/\partial a = 0$ and $\partial SS/\partial b = 0$ gives

$$\begin{aligned}\partial SS/\partial a &:= -2\sum_{i=1}^n e_i = -2\sum_{i=1}^n (y_i - a - bx_i), \\ \partial SS/\partial b &:= -2\sum_{i=1}^n x_i e_i = -2\sum_{i=1}^n x_i (y_i - a - bx_i).\end{aligned}$$

To find the minimum, we equate both these to zero:

$$\sum_{i=1}^n (y_i - a - bx_i) = 0, \quad \text{and} \quad \sum_{i=1}^n x_i (y_i - a - bx_i) = 0.$$

This gives two simultaneous linear equations in the two unknowns a, b , called the *normal equations*. Using the notation

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i,$$

dividing both sides by n and rearranging, the normal equations are

$$a + b\bar{x} = \bar{y}, \quad \text{and} \quad a\bar{x} + b\overline{x^2} = \overline{xy}.$$

Multiply the first by \bar{x} and subtract from the second:

$$b = (\overline{xy} - \bar{x}\bar{y})/(\overline{x^2} - (\bar{x})^2), \quad \text{and then} \quad a = \bar{y} - b\bar{x}.$$

We will use this bar notation systematically. We call $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$ the *sample mean*, or average, of x_1, \dots, x_n , and similarly for \bar{y} . In this book (though not all others!), the *sample variance* is defined as the average, $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$, of $(x_i - \bar{x})^2$, written s_x^2 or s_{xx} . Then using linearity of average, or ‘bar’,

$$s_x^2 = s_{xx} = \overline{(x - \bar{x})^2} = \overline{x^2 - 2x\bar{x} + \bar{x}^2} = \overline{x^2} - 2\bar{x}.\bar{x} + (\bar{x})^2 = \overline{x^2} - (\bar{x})^2.$$

Similarly, the *sample covariance* of x and y is defined as the average of $(x - \bar{x})(y - \bar{y})$, written s_{xy} . So

$$s_{xy} = \overline{(x - \bar{x})(y - \bar{y})} = \overline{xy - x.\bar{y} - \bar{x}.y + \bar{x}.\bar{y}} = \overline{(xy)} - \bar{x}.\bar{y} - \bar{x}.\bar{y} + \bar{x}.\bar{y} = \overline{(xy)} - \bar{x}.\bar{y}.$$

Thus the slope b is given by $b = s_{xy}/s_{xx}$, the ratio of the sample covariance to the sample x -variance.

Q2. With two regressors u and v and response variable y , given a sample of size n of points $(Uu_1, v_1, y_1), \dots, (u_n, v_n, y_n)$ we have to fit a least-squares *plane* – that is, choose parameters a, b, c to minimise the sum of squares

$$SS := \sum_{i=1}^n (y_i - c - au_i - bv_i)^2.$$

Taking $\partial SS/\partial c = 0$ gives

$$\sum_{i=1}^n (y_i - c - au_i - bv_i) = 0 : \quad c = \bar{y} - a\bar{u} - b\bar{v}.$$

We re-write SS as

$$SS = \sum_{i=1}^n [(y_i - \bar{y}) - a(u_i - \bar{u}) - b(v_i - \bar{v})]^2.$$

Then $\partial SS/\partial a = 0$ and $\partial SS/\partial b = 0$ give

$$\begin{aligned} \sum_{i=1}^n (u_i - \bar{u})[(y_i - \bar{y}) - a(u_i - \bar{u}) - b(v_i - \bar{v})], \\ \sum_{i=1}^n (v_i - \bar{v})[(y_i - \bar{y}) - a(u_i - \bar{u}) - b(v_i - \bar{v})]. \end{aligned}$$

Multiply out, divide by n to turn the sums into averages, and re-arrange using our earlier notation: these become

$$as_{uu} + bs_{uv} = s_{yu},$$

$$as_{uv} + bs_{vv} = s_{yv}.$$

These are the *normal equations* for a and b . The determinant is

$$s_{uu}s_{vv} - s_{uv}^2 = s_{uu}s_{vv}(1 - r_{uv}^2)$$

(as $r_{uv} := s_{uv}/(s_u \cdot s_v)$), $\neq 0$ iff $r_{uv} \neq \pm 1$, i.e., iff the (u_i, v_i) are not collinear, and this is the condition for the normal equations to have a unique solution.

NHB