

Tutorial # 11

N 22.5. 1-5

Regression analysis:

We measure a random variable Y at fixed pts $X = [x_1, \dots, x_n]$.

Suppose we have just 1 sample:

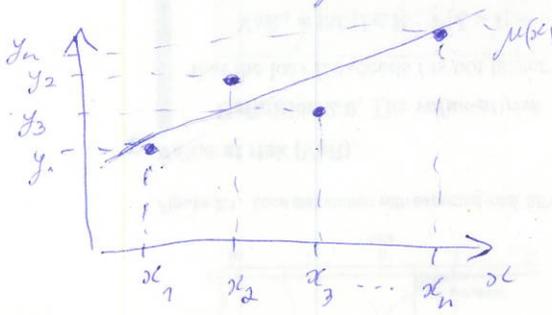
values $[y_1, \dots, y_n]$ measured at $[x_1, \dots, x_n]$ pts.

In linear regression, we assume mean value of Y

μ depends on x as $\mu(x) = a + bx$

regression line regression coefficient

For 1 sample, we have:



roughly speaking, we want to pass a line exactly "in between" all measured values $[y_1, \dots, y_n]$.

That is, we want to find a, b s.t. the deviation of y_1, \dots, y_n from regression line $\mu(x)$ is minimal.

$S(a, b) = \sum_{j=1}^n (y_j - a - bx_j)^2$ should be minimized.

$\frac{\partial S}{\partial a} = 0 \Rightarrow -2 \cdot \sum_{j=1}^n (y_j - a - bx_j) = 0 \Rightarrow a = \frac{\sum_{j=1}^n y_j}{n} - b \cdot \frac{\sum_{j=1}^n x_j}{n}$

$\frac{\partial S}{\partial b} = 0 \Rightarrow -2 \cdot \sum_{j=1}^n x_j (y_j - a - bx_j) = 0 \Rightarrow b = \frac{\sum_{j=1}^n x_j y_j - a \sum_{j=1}^n x_j}{\sum_{j=1}^n x_j^2}$

$$\begin{cases} a = \bar{y} - b\bar{x} \\ b = \end{cases}$$

$$b = \frac{\sum_{j=1}^n x_j y_j - a n \bar{x}}{\sum_{j=1}^n x_j^2} = \frac{\sum_{j=1}^n x_j y_j - n \bar{x} \bar{y}}{\sum_{j=1}^n x_j^2} + \frac{b n \bar{x}^2}{\sum_{j=1}^n x_j^2}$$

$$\Rightarrow b = \frac{\sum_{j=1}^n x_j y_j - n \bar{x} \bar{y}}{\sum_{j=1}^n x_j^2 - n \bar{x}^2}, \quad a = \bar{y} - b \bar{x}$$

It's conventional to rewrite this in terms

of "variance / covariance" (even though x is not a random variable) :

$$S_{xx}^2 := \frac{1}{n-1} \cdot \sum_{j=1}^n (x_j - \bar{x})^2$$

$$S_{xy} := \frac{1}{n-1} \cdot \sum_{j=1}^n (x_j - \bar{x})(y_j - \bar{y})$$

$$S_{xx}^2 = \frac{1}{n-1} \cdot \sum_{j=1}^n x_j^2 - \frac{2\bar{x}}{n-1} \cdot \sum_{j=1}^n x_j + \frac{n \cdot \bar{x}^2}{n-1} = \frac{1}{n-1} \cdot \sum_{j=1}^n x_j^2 - \frac{n \bar{x}^2}{n-1}$$

$$S_{xy} = \frac{1}{n-1} \left[\sum_{j=1}^n x_j y_j - \bar{x} \cdot \sum_{j=1}^n y_j - \bar{y} \cdot \sum_{j=1}^n x_j + n \bar{x} \bar{y} \right] = \frac{1}{n-1} \cdot \sum_{j=1}^n x_j y_j - n \cdot \frac{\bar{x} \bar{y}}{n-1}$$

$$\Rightarrow \begin{cases} \sum_{j=1}^n x_j^2 - n \bar{x}^2 = (n-1) \cdot S_{xx}^2 \\ \sum_{j=1}^n x_j y_j - n \bar{x} \bar{y} = (n-1) \cdot S_{xy} \end{cases}$$

$$\Rightarrow b = \frac{S_{xy}}{S_{xx}^2}$$

Also, introduce $S_y^2 := \frac{1}{n-1} \cdot \sum_{j=1}^n (y_j - \bar{y})^2$.

Similarly as for S_x^2 , we have

$$\sum_{j=1}^n y_j^2 - n\bar{y}^2 = (n-1)S_y^2$$

- don't need, in fact

Then, note that we can write

$$\begin{aligned} S(a, b) &= \sum_{j=1}^n (y_j - a - bx_j)^2 = \sum_{j=1}^n [y_j - \bar{y} - b(x_j - \bar{x})]^2 = \\ &= \sum_{j=1}^n (y_j - \bar{y})^2 + b^2 \sum_{j=1}^n (x_j - \bar{x})^2 - 2b \sum_{j=1}^n (y_j - \bar{y})(x_j - \bar{x}) = \\ &= (n-1) \cdot \left[S_y^2 + b^2 S_x^2 - 2b S_{xy} \right] = (n-1) \cdot \left(S_y^2 - \frac{S_{xy}^2}{S_x^2} \right) \geq 0 \end{aligned}$$

$$\Leftrightarrow S_y^2 \geq \frac{S_{xy}^2}{S_x^2} \Leftrightarrow \frac{|S_{xy}|}{|S_x| |S_y|} \leq 1$$

moreover, $S_y^2 = \frac{S_{xy}^2}{S_x^2}$

$$\Leftrightarrow S(a, b) = 0,$$

i.e. all the sample pts lie on one line

$$|S_{xy}| = |S_x| |S_y|$$

Correlation:

Now, let both X and Y be random variables.

Similarly to b , define correlation coefficient of the sample:

$$r = \frac{S_{xy}}{S_x S_y} \quad (s_x > 0, s_y > 0)$$

This gives an estimate for correlation function of a population:

$$\rho = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

covariance

$$\sigma_{xy} = E[(X - \mu_x)(Y - \mu_y)]$$

$$\sigma_x = E[(X - \mu_x)^2]$$

$$\sigma_y = E[(Y - \mu_y)^2]$$

variances

(this measures linear (!) dependence of X and Y)

Note that S_{xy}, S_x, S_y introduced above are exactly unbiased estimates of $\sigma_{xy}, \sigma_x, \sigma_y$, respectively.

$$0 \leq |\rho| \leq 1 \quad \Leftrightarrow \quad 0 \leq |\tau| \leq 1$$

exact linear dependence: deviation is zero
no linear dependence: $(X - \mu_x)(Y - \mu_y)$ has positive & negative value equally likely

Compare:

In case of regression, $S_{xy} = 0 \Leftrightarrow b = 0$ implies Y being constant (i.e. independent) in X .