

Least time.  
Sample variance  
See 11.5 Least square linear regression

Today.

Some properties of  $r_{X,Y} = \frac{\mathbb{E} X^* Y^*}{\sigma_X \sigma_Y} = \frac{\mathbb{E} D_Y D_X}{\sigma_X \sigma_Y}$

① Symmetric about X & Y.

$r_{X,Y} = r_{Y,X} = \underbrace{\quad}_{\text{good as long as no ambiguity.}}$

Sec 11.4

② Bounds

by Cauchy-Schwarz inequality.

$\mathbb{E} X Y \leq \sqrt{\mathbb{E} X^2 \mathbb{E} Y^2}$

$r = \frac{\mathbb{E} X^* Y^*}{\sigma_X \sigma_Y} \leq \frac{\sqrt{\mathbb{E} X^{*2} \mathbb{E} Y^{*2}}}{\sigma_X \sigma_Y} = 1$

lower bound for r  
= - upper bound for (-r)

$-r \geq \frac{\mathbb{E} [(-X^*) \cdot Y^*]}{\sqrt{\mathbb{E} (-X^*)^2 \mathbb{E} Y^{*2}}} = -1$

$= 1 \Rightarrow r \geq -1$

Proof w/o using C-S:

Consider  $\mathbb{E} (X^* - Y^*)^2 \geq 0$

$\mathbb{E} X^{*2} - 2 \mathbb{E} X^* Y^* + \mathbb{E} Y^{*2} \geq 0$

$-2 \mathbb{E} X^* Y^* \leq \mathbb{E} X^{*2} + \mathbb{E} Y^{*2} \geq 2$

$r \geq \mathbb{E} X^* Y^* \leq 1$

For lower bound. —

$r \geq -1$

bounds for r : -1, 1

When  $r=1$ ?  $\mathbb{E} (X^* - Y^*)^2 = 0$

$\Rightarrow X^* - Y^* = 0$

$\Rightarrow \frac{X - \mu_X}{\sigma_X} = \frac{Y - \mu_Y}{\sigma_Y}$

$\Rightarrow X$  is a linear function of  $Y$

$X = aY + b$

$r=1 \Leftrightarrow \begin{matrix} a > 0 \\ a < 0 \end{matrix}$

Sec 11.5 The error in regression

Error in regression

$D_Y = Y - \hat{Y}$

least square linear regression estimator of Y based on X

Residual, the left-over part that cannot be estimated

$Y = \hat{Y} + D$

$\hat{Y} = \hat{a}X + \hat{b}$

$\hat{a} = r \frac{\sigma_Y}{\sigma_X}, \hat{b} = \mu_Y - \hat{a} \mu_X$

$= \hat{a}X + \mu_Y - \hat{a} \mu_X$

$D = Y - (\hat{a}X + \mu_Y - \hat{a} \mu_X)$

$= (Y - \mu_Y) - \hat{a} (X - \mu_X)$

$= D_Y - \hat{a} D_X$  — no intercept.

$\mathbb{E} D = \underbrace{\mathbb{E} D_Y}_{=0} - \hat{a} \underbrace{\mathbb{E} D_X}_{=0} = 0$

This means, average residual is always 0, no matter what the shape of scatter diagram is } describing from Data 8

Or, in the prob. word,

no matter what the joint distr. of (X,Y) is



$\text{Var}(D) = \mathbb{E} D^2 = \text{MSE}$

Recall  $\mathbb{E} D^2 = \mathbb{E} (D_Y - \hat{a} D_X)^2$

$= \mathbb{E} D_Y^2 - 2 \hat{a} \mathbb{E} D_Y D_X + \hat{a}^2 \mathbb{E} D_X^2$

$= \sigma_Y^2 - 2 \hat{a} r \sigma_X \sigma_Y + \hat{a}^2 \sigma_X^2$

$= \sigma_Y^2 - 2 r \frac{\sigma_Y}{\sigma_X} \sqrt{\sigma_X^2 \sigma_Y^2} + r^2 \frac{\sigma_Y^2}{\sigma_X^2} \sigma_X^2$

$= \sigma_Y^2 - 2 r^2 \sigma_Y^2 + r^2 \sigma_Y^2$

$= (1 - r^2) \sigma_Y^2$

$\text{SD}(D) = \sqrt{1 - r^2} \sigma_Y$  — same formula as in Data 8

Since  $\mathbb{E} D = 0$ , when  $\text{SD}(D)$  is small, by Chebyshev's inequality, D is with high prob. close to 0

i.e.  $\hat{Y} \approx Y$  with high prob.  $\Rightarrow$  good estimation.

Another extreme case,  $r=0$   $\text{SD}(D) = \sigma_Y$

$\hat{a} = r \frac{\sigma_Y}{\sigma_X} = 0 \Rightarrow$  the best linear estimator of Y using X

= the best constant estimator

$\Rightarrow$  Do not need  $X$  /  $X$  provides no extra info.

when we estimate Y linearly

To summarize, |r| quantifies the amount of linear association between X and Y

Note: even when  $r=0$ , it is totally possible that X and Y have a non-linear association



The residual D, and X

Data 8: residual plot — flat

(Imagine that you are doing another linear regression of D) based on X

Slope =  $\hat{a}_{D,X} = r_{D,X} \frac{\sigma_D}{\sigma_X}$

WTS slope  $\Rightarrow$

$r_{D,X} = \frac{\text{Cov}(D, X)}{\sigma_D \sigma_X} = 0$

$\text{Cov}(D, X) = \mathbb{E} (D - \mathbb{E} D) (X - \mathbb{E} X) = \mathbb{E} D D_X$

$= \mathbb{E} [D_X (D_Y - \hat{a} D_X)]$

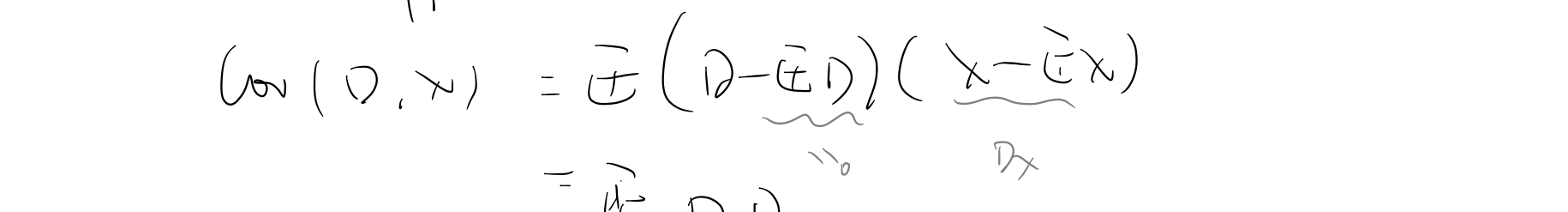
$= \mathbb{E} D_X D_Y - \hat{a} \mathbb{E} D_X^2$

$= r \sigma_X \sigma_Y - r \frac{\sigma_Y}{\sigma_X} \cdot \sigma_X^2 = 0$

intercept =  $\mu_D - \hat{a}_{D,X} \mu_X = \mathbb{E} D = 0 \Rightarrow 0$

Sec 12.1 Simple linear regression model

machine / black box



What we know: input & output

Wish to learn about the machine.

i.e. linear params.

Formally, n (input, output)  $(x_i, y_i) \quad i=1, 2, \dots, n$

$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad i=1, 2, \dots, n$

$\beta_0, \beta_1$  are unknown params that we wish to estimate

$x_i$ 's are considered known, fixed value

$\epsilon_i \sim N(0, \sigma^2)$  (same var. for all i)

$\Rightarrow Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$  indep

We wish to give an estimator of Y as

$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

Error  $= Y_i - \hat{Y}_i$

Target: find  $(\hat{\beta}_0, \hat{\beta}_1)$  to minimizing MSE

$= \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$

$Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$  indep

"The randomness of  $Y_i$  comes only from the noise  $\epsilon_i$

Therefore, the average output is

$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i \sim N(\beta_0 + \beta_1 \bar{X}, \frac{\sigma^2}{n})$

As linear combination of indep normals

$\mathbb{E} \bar{Y} = \frac{1}{n} \sum_{i=1}^n \mathbb{E} Y_i = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i)$

$= \frac{1}{n} \left[ n \beta_0 + \beta_1 \sum_{i=1}^n x_i \right]$

$= \beta_0 + \beta_1 \bar{X}$

"the same linear transform of the average input"

$\text{Var}(\bar{Y}) = \text{Var}(\text{Constant} + \frac{1}{n} \sum_{i=1}^n \epsilon_i)$

$= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(\epsilon_i)$

$= \frac{1}{n^2} \sum_{i=1}^n \sigma^2$

$= \frac{1}{n^2} \cdot n \sigma^2 = \frac{\sigma^2}{n}$

Define them before using them

$\frac{X - \mathbb{E} X}{\text{SD}(X)} = \frac{X - \mu_X}{\sigma_X} = \frac{D_X}{\sigma_X}$

$X^*$  is X in standard units

$\Rightarrow \mathbb{E} X^* = 0, \text{SD}(X^*) = 1$

$\mathbb{E} X^* = \mathbb{E} \frac{X - \mathbb{E} X}{\text{SD}(X)}$

$= \frac{1}{\text{SD}(X)} (\mathbb{E} X - \mathbb{E} X) = 0$

$\text{SD}(X^*) = \text{SD}\left(\frac{X - \mathbb{E} X}{\text{SD}(X)}\right)$

$= \frac{1}{\text{SD}(X)} \text{SD}(X) = 1$

$\text{Var}(X^*) = \mathbb{E} X^{*2} - (\mathbb{E} X^*)^2 = 1$

$\Rightarrow \mathbb{E} X^{*2} = \text{Var}(X^*) = 1^2 = 1$

Same for  $\mathbb{E} Y^{*2} = 1$