

Recall:  $M_{Y|X} =$  Regression of  $Y$  on  $X$

$E(Y|X)$  = conditional expectation.

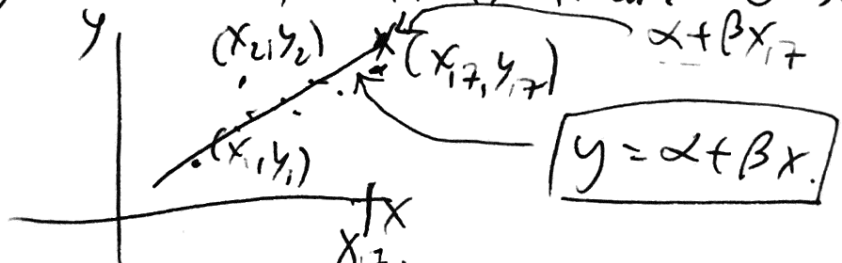
If linear,  $M_{Y|X} = \alpha + \beta X$ ; then.

$$M_{Y|X} = \nu + \frac{c}{\sigma^2} (x - \mu), \text{ where}$$

$$\nu = E(Y), \mu = E(X), \sigma^2 = \text{Var}(X), c = \text{Cov}(X, Y).$$

i.e.  $M_{Y|X} = \underbrace{\nu - \mu \frac{c}{\sigma^2}}_{\alpha} + \underbrace{\frac{c}{\sigma^2}}_{\beta} X$ .

Assume regression of  $Y$  on  $X$  is linear. Observe  $(x_i, y_i)_{i=1}^n$



What should we optimize? What are the errors?

Each error looks like:  $(\alpha + \beta x_i) - y_i$

Legendre: Optimize  $\sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 = Q$ .

To find minimum of  $Q$ , solve  $\frac{\partial Q}{\partial \alpha} = 0 = \frac{\partial Q}{\partial \beta}$ .

11

Will give "least square error" estimator,  $\hat{\alpha}, \hat{\beta}$ .

$$\frac{\partial Q}{\partial \alpha} = \sum_{i=1}^n \frac{\partial}{\partial \alpha} (y_i - (\alpha + \beta x_i))^2 = 0.$$

$$\Leftrightarrow n\bar{y} - n\alpha - \beta n\bar{x} = 0$$

$$\Leftrightarrow \alpha + \beta\bar{x} = \bar{y} \quad \Leftrightarrow \alpha = \bar{y} - \beta\bar{x}$$

$$0 = \frac{\partial Q}{\partial \beta} = \sum_{i=1}^n \frac{\partial}{\partial \beta} (y_i - (\alpha + \beta x_i))^2 (x_i) = 0;$$

$$\sum_{i=1}^n x_i y_i = \alpha n\bar{x} + \beta \sum_{i=1}^n x_i^2$$

$$\text{Let } S_{xy} := \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}.$$

$$S_{xx} := \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2.$$

$$\rightarrow S_{xy} + n\bar{x}\bar{y} = \alpha n\bar{x} + \beta (S_{xx} + n\bar{x}^2).$$

$$-n\bar{x}\bar{y} = -\alpha n\bar{x} - n\beta\bar{x}^2$$

$$S_{xy} = \beta S_{xx},$$

$$\hat{\beta} = \frac{S_{xy}}{S_{xx}}, \quad \hat{\alpha} = \bar{y} - \frac{S_{xy}}{S_{xx}} \bar{x}.$$

How good are  $\hat{\alpha}$  &  $\hat{\beta}$  as estimators for  $\alpha$  &  $\beta$ ?

Extra assumption: normal (linear) regression analysis:

$f_Y(y|X=x)$  = conditional density,

$$M_{Y|X} = \alpha + \beta x.$$

$$E(Y|X)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-(\alpha+\beta x))^2}{2\sigma^2}}$$

With this extra assumption, can

compute max likelihood estimators for  $\hat{\alpha}$ ,  $\hat{\beta}$ ,  $\hat{\sigma}$ .

Observed  $(x_i, y_i)_{i=1}^n$

$$L(y_1, \dots, y_n | x_1, \dots, x_n, \alpha, \beta, \sigma) = \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right)^n e^{-\frac{(y_i - (\alpha + \beta x_i))^2}{2\sigma^2}}$$

$$\frac{\partial L}{\partial \alpha} = 0 = \frac{\partial \ln L}{\partial \beta} = \frac{\partial \ln L}{\partial \sigma} \quad \hookrightarrow \text{same with } \log L.$$

$$\log L = \underbrace{-\frac{n}{2} \log 2\pi}_{\text{constant}} - \underbrace{n \log \sigma}_{\text{constant}} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2$$

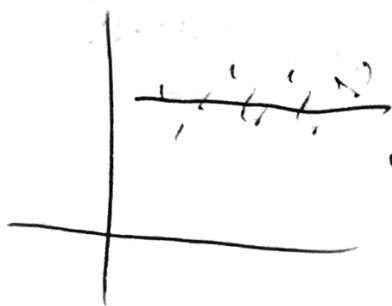
$$\frac{\partial \log L}{\partial \alpha} = 0 = \frac{\partial \log L}{\partial \beta} \Rightarrow \hat{\alpha} \text{ \& \ } \hat{\beta} \text{ as before.}$$

$$0 = \frac{\partial \log L}{\partial \sigma} = \frac{-n\sigma^2}{\sigma^3} + \frac{1}{\sigma^3} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2 = 0.$$

$$\hat{\sigma}^2 = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta}x_i))^2}$$

(exercise)  $= \sqrt{\frac{1}{n} (S_{yy} - \hat{\beta} S_{xy})}$

Let's test (under normality assumption) ~~the~~ hypotheses on parameters, I.e.  $\hat{\beta} = 0$ ?



$\frac{S_{xy}}{S_{xx}}$  ← correlation

Null hypothesis (if we want to show that the variables are

correlated).

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

$y_i =$  observed values  
 $Y_i =$  RV

Let  $\hat{\beta} = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_{xx}} \right) (Y_i - \bar{Y}) =$  RV.

$= \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_{xx}} \right) Y_i$  ← normal

Normal!

Need: mean & variance  $\hat{\beta}$ .

$$\begin{aligned} E(\hat{\beta}) &= E(\hat{\beta} | X_i) = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_{xx}} \right) \underbrace{E(Y_i | x_i)}_{(\alpha + \beta x_i)} \\ &= \beta \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_{xx}} \right) x_i = \frac{\beta}{S_{xx}} \left[ \sum_{i=1}^n x_i^2 - \bar{x} n \bar{x} \right] \\ &= \beta. \end{aligned}$$

---

Next:  $\text{Var}(\hat{\beta}) = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_{xx}} \right)^2 \underbrace{\text{Var}(Y_i)}_{\sigma^2}$

$$= \frac{\sigma^2}{S_{xx}^2} \cdot \underbrace{\sum_{i=1}^n (x_i - \bar{x})^2}_{S_{xx}} = \boxed{\frac{\sigma^2}{S_{xx}}}$$

---

Recap:  $\hat{\beta} = \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{S_{xx}} \right) Y_i \quad \Rightarrow$

normal random variable with mean  $E(\hat{\beta}) = \beta$ .

&  $\text{Var}(\hat{\beta}) = \frac{\sigma^2}{S_{xx}}$ .

Fact:  $\frac{n \hat{\sigma}^2}{\sigma^2} = \frac{1}{\sigma^2} (S_{YY} - \hat{\beta} S_{XY}) \sim \chi^2_{n-2}$  deg.

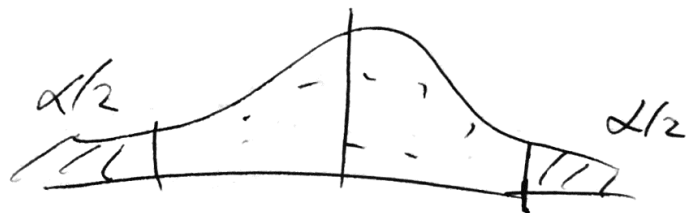
&  $\frac{n \hat{\sigma}^2}{\sigma^2}$  is indep of  $\hat{\beta}$ .

(5)

$$Z = T = \frac{\hat{\beta} - \beta}{\sqrt{\frac{\hat{\sigma}^2}{S_{xx}}}} \cdot \sqrt{\frac{\hat{\sigma}^2 (n-2)}{n \hat{\sigma}^2}}$$

Then:  $(\hat{\beta} - \beta) \cdot \sqrt{\frac{S_{xx}(n-2)}{\hat{\sigma}^2 n}}$  is a T-distributed RV with  $n-2$  degrees of freedom.

Q: Construct a  $(1-\alpha)100\%$  confidence interval for  $\beta$ :



$$P\left(-t_{\alpha/2, n-2} < (\hat{\beta} - \beta) \sqrt{\frac{(n-2) S_{xx}}{n \hat{\sigma}^2}} < t_{\alpha/2, n-2}\right) = 1 - \alpha.$$

$$\hat{\beta} - \underbrace{\quad} < \beta < \hat{\beta} + \underbrace{t_{\alpha/2, n-2} \sqrt{\frac{n \hat{\sigma}^2}{(n-2) S_{xx}}}}_{\quad}$$