

Exam problem: profit/overbooking.

$$e_i = y_i - (\alpha + \beta x_i)$$

Least squares estimates $\hat{\alpha}^*$ and $\hat{\beta}^*$ minimize $\sum_{i=1}^n e_i^2$

Why e_i^2 ? it is easier than $|e_i|$.

Formula: Theorem 14.2

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

$$\text{where } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

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

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

Proof:

$$\text{let } q = \sum_{i=1}^n (y_i - (\hat{\alpha} + \hat{\beta} x_i))^2$$

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

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

solving system for $\hat{\alpha}, \hat{\beta}$, we get desired result.

Ex: hours studied vs test score:

$$\bar{x} = 10 \quad \bar{y} = 56.4$$

$$S_{xx} = 376, \quad S_{xy} = 1305$$

$$\hat{\beta}^* = 3.471 \quad \hat{\alpha}^* = 21.69$$

so least-squares line is $y = 21.69 + 3.471x$

if you study 14 hours you should get $21.69 + 3.471 \cdot 14$
 $= 70.284$

Remark: x could be multidimensional. i.e. linear function of several variables.

§14.4: Normal Regression Analysis.

recall that $\mu_{y|x} = E(y|x) = \int_{-\infty}^{\infty} y w(y|x) dy$

In particular, we assume

$$w(y|x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2} (y - (\alpha + \beta x))^2\right]$$

(normal density with $\mu = \alpha + \beta x$ and variance σ^2)

Now given $\{(x_i, y_i) : i=1, \dots, n\}$ we'd like to

① the estimation of σ, α, β .

② Test of hypotheses concerning these

③ predictions based on estimate $\mu_{y|x} = \hat{\alpha} + \hat{\beta}x$

① MLE of α, β, σ

① MLE of α, β, σ

likelihood $\prod_{i=1}^n w(y_i | x_i) \cdot m(x_i)$

← marginal density of x
doesn't contain α, β, σ .

$f(x_i, y_i)$

So for simplicity let $L = \prod_{i=1}^n w(y_i | x_i)$

$$\log L = \sum_{i=1}^n \log(w(y_i | x_i))$$

$$= \sum_{i=1}^n \left(-\log(\sigma) - \frac{1}{2} (\log(2\pi)) - \frac{1}{2\sigma^2} (y_i - (\alpha + \beta x_i))^2 \right)$$

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

Notice that the MLE estimates of α, β are the $\arg \max_{\alpha, \beta}$ (this)

$$= \arg \min_{\alpha, \beta} \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2$$

↓ ↓

least squares estimate

$$\begin{matrix} \hat{\alpha}^* & \hat{\beta}^* \\ \downarrow & \downarrow \\ \bar{y} - \hat{\beta}^* \bar{x} & \frac{S_{xy}}{S_{xx}} \end{matrix}$$

So MLE = LSE

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

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

$$= \frac{1}{n} \sum_{i=1}^n (y_i - (\bar{y} - \hat{\beta} \bar{x} + \hat{\beta} x_i))^2$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{i=1}^n (y_i - (\bar{y} - \hat{\beta}\bar{x} + \hat{\beta}x_i))^2 \\
&= \frac{1}{n} \left[\sum_{i=1}^n (y_i - \bar{y})^2 - 2\hat{\beta} \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + \hat{\beta}^2 \sum_{i=1}^n (x_i - \bar{x})^2 \right] \\
&\quad \quad \quad \uparrow \quad \quad \quad \uparrow \quad \quad \quad \uparrow \\
&\quad \quad \quad s_{yy} \quad \quad \quad s_{xy} \quad \quad \quad s_{xx} \\
&= \frac{1}{n} [s_{yy} - 2\hat{\beta}s_{xy} + \hat{\beta}s_{xx}] \\
&= \frac{1}{n} [s_{yy} - \hat{\beta}s_{xy}]
\end{aligned}$$

$$\Rightarrow \hat{\sigma} = \sqrt{\frac{1}{n}(s_{yy} - \hat{\beta}s_{xy})}$$

Remarks:

1: for the Normal Regression,

$$w(y|x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(y - (\alpha + \beta x))^2\right]$$

the regression on $\mu_{y|x} = \alpha + \beta x$ is linear.

hence by 14.1 $\alpha = \mu_2 - \rho \frac{\sigma_2}{\sigma_1} \mu_1$, $\beta = \rho \frac{\sigma_2}{\sigma_1}$

Compare w/ MLE.

$$\alpha = \mu_2 - \rho \frac{\sigma_2}{\sigma_1} \mu_1$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad \checkmark$$

$$\rho = \rho \frac{\sigma_2}{\sigma_1} = \frac{\text{Cov}(x, y)}{\sigma_1^2}$$

$$\hat{\beta} = \frac{s_{xy}}{s_{xx}} \quad \checkmark$$