

## Approximation of mean and variance for functions of random variables.

Example. Measure voltage,  $U$ , and current,  $I$ , and want to calculate the resistance  $R = \frac{U}{I}$ . What about  $E[R]$  and  $\text{Var}[R]$ .

Let  $E[U] = \mu_1$  and  $E[I] = \mu_2$

From a Taylor-expansion of  $f(x,y) = \frac{x}{y}$  around  $(\mu_1, \mu_2)$  we get,

$$\begin{aligned} f(x,y) &= f(\mu_1, \mu_2) + \frac{\partial f}{\partial x} \Big|_{(\mu_1, \mu_2)} (x - \mu_1) + \frac{\partial f}{\partial y} \Big|_{(\mu_1, \mu_2)} (y - \mu_2) + \varepsilon \\ &= \frac{\mu_1}{\mu_2} + \frac{1}{\mu_2} (x - \mu_1) - \frac{\mu_1}{\mu_2^2} (y - \mu_2) + \varepsilon \end{aligned}$$

i.e. if  $\varepsilon$  is small,  $R$  can be approximated by

$$R = \frac{U}{I} \approx \frac{\mu_1}{\mu_2} + \frac{1}{\mu_2} (U - \mu_1) - \frac{\mu_1}{\mu_2^2} (I - \mu_2)$$

such that

$$E[R] \approx \frac{\mu_1}{\mu_2} + \frac{1}{\mu_2} E[U - \mu_1] - \frac{\mu_1}{\mu_2^2} E[I - \mu_2] = \frac{\mu_1}{\mu_2}$$

and

$$\text{Var}[R] \approx \left[ \frac{1}{\mu_2} \right]^2 \text{Var}[U] + \frac{\mu_1^2}{\mu_2^4} \text{Var}[I] - \frac{2\mu_1}{\mu_2^3} \text{Cov}[U, I]$$

If the uncertainty in  $U$  and  $I$  ~~are~~ are measurement errors it is reasonable that  $\text{Cov}[U, I] = 0$

In general if  $Y = f(x_1, \dots, x_n)$ , we have,

$$f(x_1, \dots, x_n) = f(\mu_1, \dots, \mu_n) + \sum_{j=1}^n \frac{\partial f}{\partial x_j} \Big|_{(\mu_1, \dots, \mu_n)} (x_j - \mu_j) + \varepsilon$$

## Approximation of mean and variance for functions of random variables.

Example. Measure voltage,  $U$ , and current,  $I$ , and want to calculate the resistance  $R = \frac{U}{I}$ . What about  $E[R]$  and  $\text{Var}[R]$ .

$$\text{Let } E[U] = \mu_1 \text{ and } E[I] = \mu_2$$

From a Taylor-expansion of  $f(x, y) = \frac{x}{y}$  around  $(\mu_1, \mu_2)$  we get,

$$\begin{aligned} f(x, y) &= f(\mu_1, \mu_2) + \frac{\partial f}{\partial x} \Big|_{(\mu_1, \mu_2)} (x - \mu_1) + \frac{\partial f}{\partial y} \Big|_{(\mu_1, \mu_2)} (y - \mu_2) + \varepsilon \\ &= \frac{\mu_1}{\mu_2} + \frac{1}{\mu_2} (x - \mu_1) - \frac{\mu_1}{\mu_2^2} (y - \mu_2) + \varepsilon \end{aligned}$$

i.e. if  $\varepsilon$  is small,  $R$  can be approximated by

$$R = \frac{U}{I} \approx \frac{\mu_1}{\mu_2} + \frac{1}{\mu_2} (U - \mu_1) - \frac{\mu_1}{\mu_2^2} (I - \mu_2)$$

such that

$$E[R] \approx \frac{\mu_1}{\mu_2} + \frac{1}{\mu_2} E[U - \mu_1] - \frac{\mu_1}{\mu_2^2} E[I - \mu_2] = \frac{\mu_1}{\mu_2}$$

$$\text{and } \text{Var}[R] \approx \left[ \frac{1}{\mu_2} \right]^2 \text{Var}[U] + \frac{\mu_1^2}{\mu_2^4} \text{Var}[I] - \frac{2\mu_1}{\mu_2^3} \text{Cov}[U, I]$$

If the uncertainty in  $U$  and  $I$  ~~are~~ are measurement errors it is reasonable that  $\text{Cov}[U, I] = 0$

In general if  $Y = f(x_1, \dots, x_m)$ , we have:

$$f(x_1, \dots, x_m) = f(\mu_1, \dots, \mu_m) + \sum_{j=1}^m \frac{\partial f}{\partial x_j} \Big|_{(\mu_1, \dots, \mu_m)} (x_j - \mu_j) + \varepsilon$$

and if higher order terms are small, we have

$$E[Y] \approx f(u_1, \dots, u_n)$$

$$\text{Var}[Y] \approx \sum_{i=1}^n \frac{\partial f}{\partial x_i}(u_1, \dots, u_n)^2 \text{Var}[X_i] + 2 \sum_{j < k} \frac{\partial f}{\partial x_j}(u_1, \dots, u_n) \frac{\partial f}{\partial x_k}(u_1, \dots, u_n) \text{Cov}(X_j, X_k)$$

Example  $U \sim N(100, 2^2)$   $R = \frac{U}{I}$   
 $I \sim N(20, 1)$

We get,  $E[R] \approx 5$  and  $\text{Var}[R] \approx \left(\frac{1}{20}\right)^2 \cdot 4 + \left(\frac{100}{400}\right)^2 \cdot 1$

$$= \frac{4}{400} + \frac{1}{16} = \frac{1}{100} + \frac{1}{16} = 0.077$$

$$\begin{cases} U \sim N(100, 5^2) \\ I \sim N(20, 1^2) \end{cases} \Rightarrow E[R] \approx 5 \\ \text{Var}[R] \approx \left(\frac{1}{20}\right)^2 \cdot 25 + \left(\frac{100}{400}\right)^2 \cdot 4 \approx 0.31$$

Purpose: Find relationships between one response variable and one or more explanatory variables. Regression analysis is a modelling tool.

## 11.2. Simple linear regression (linear in the coefficients)

$$\text{Model: } Y_i = \alpha + \beta X_i + \varepsilon_i \quad \left. \begin{array}{l} \text{independent} \\ E[\varepsilon_i] = 0, \quad \text{Var}[\varepsilon_i] = \sigma^2 \end{array} \right\}$$

response                      regression variable  
 dependent variable            independent variable  
 explanatory variable

$X_i$  may be a transformation of other variables.

i.e.  $X_i = \ln(\theta_i)$ ,  $X_i = \sin(\theta_i)$  and so on.

Observation:  $(y_1, x_1), (y_2, x_2), \dots, (y_m, x_m)$  which ~~we~~ have to fulfill.

$$\begin{aligned}
 y_1 &= \alpha + \beta x_1 + \varepsilon_1 \\
 y_2 &= \alpha + \beta x_2 + \varepsilon_2 \quad \text{or} \\
 &\vdots \\
 y_m &= \alpha + \beta x_m + \varepsilon_m
 \end{aligned}
 \quad \left[ \begin{array}{c} y_1 \\ y_2 \\ \vdots \\ y_m \end{array} \right] = \left[ \begin{array}{cc} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{array} \right] \left[ \begin{array}{c} \alpha \\ \beta \end{array} \right] + \left[ \begin{array}{c} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_m \end{array} \right]$$

We need to find estimates,  $(\hat{\alpha}, \hat{\beta})$ , for  $(\alpha, \beta)$ .

### Method of Least Squares

Find the  $\alpha$  and  $\beta$  that minimize

$$\alpha = \sum_{i=1}^m \varepsilon_i^2 = \sum_{i=1}^m (y_i - \alpha - \beta x_i)^2$$

$$\frac{\partial Q}{\partial a} = 0 \iff -2 \sum_{i=1}^n (y_i - a - b x_i) = 0$$

$$\frac{\partial Q}{\partial b} = 0 \iff -2 \sum_{i=1}^n (y_i - a - b x_i) x_i = 0$$

$\frac{\partial Q}{\partial b}$

This gives the 2 normal equations

$$\sum_{i=1}^m a + \sum_{i=1}^m b x_i = \sum_{i=1}^m y_i \quad (1)$$

$$\sum_{i=1}^m a x_i + \sum_{i=1}^m b x_i^2 = \sum_{i=1}^m x_i y_i \quad (2)$$

From (1) we get:

$$ma + b \sum_{i=1}^m x_i = \sum_{i=1}^m y_i \text{ which implies } a = \bar{y} - b \bar{x}$$

and substitution of a with  $\bar{y} - b \bar{x}$  in (2) gives

$$(\bar{y} - b \bar{x}) \sum_{i=1}^m x_i + b \sum_{i=1}^m x_i^2 = \sum_{i=1}^m x_i y_i$$

$$\Rightarrow b \left( \sum_{i=1}^m x_i^2 - m \bar{x}^2 \right) = \sum_{i=1}^m x_i (y_i - \bar{y})$$

$$\Rightarrow b = \frac{\sum_{i=1}^m x_i (y_i - \bar{y})}{\sum_{i=1}^m x_i^2 - m \bar{x}^2} = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2} = \frac{\sum_{i=1}^m (x_i - \bar{x}) y_i}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

$$\text{Here is used: } \sum_{i=1}^m \bar{x} (y_i - \bar{y}) = \bar{x} \sum_{i=1}^m (y_i - \bar{y}) = 0$$

$$\sum_{i=1}^m (x_i - \bar{x})^2 = \sum_{i=1}^m (x_i^2 - 2x_i \bar{x} + \bar{x}^2) = \sum_{i=1}^m x_i^2 - 2m \bar{x}^2 + m \bar{x}^2$$

$$= \sum_{i=1}^m x_i^2 - m \bar{x}^2$$

$$\sum_{i=1}^m \bar{y} (x_i - \bar{x}) = \bar{y} \sum_{i=1}^m (x_i - \bar{x}) = 0$$

The corresponding estimators are:  $\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}$

and  $\hat{a} = \bar{y} - \hat{\beta} \bar{x}$ .

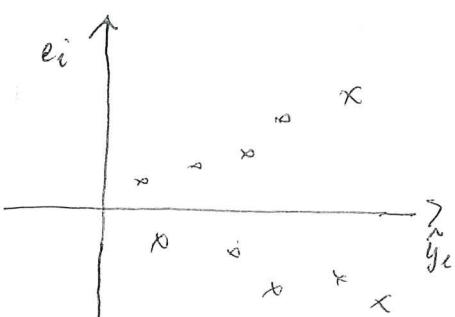
The estimated expected regression line (estimated model) is:

$$\hat{y} = a + b x$$

$$\text{Residuals: } e_i = y_i - \hat{y}_i = y_i - a - b x_i$$

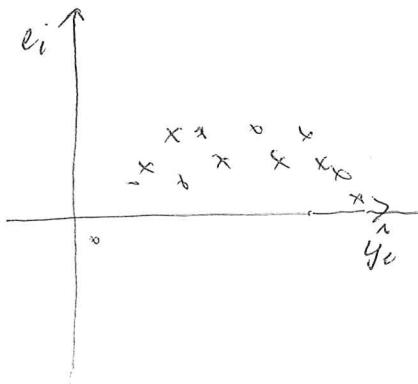
These should be plotted against  $\hat{y}_i$ , against  $x_i$  and eventually other regression variables. Normal-plot should be used to check for normal distribution.

Pattern in the residuals indicates that the model does not fit.



Variance not constant.

Transform  $y$

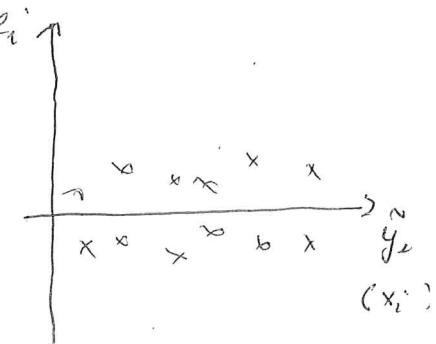


Should consider

second order

terms

$$\therefore E[y_i] = a + \beta_1 x_i + \beta_2 x_i^2$$



No pattern

(ok.)