

Supplementary lecture note

1. Modeling relationships between the rv's X and Y

We saw in the lecture that when modeling a relationship between two random variables (rv's), X and Y (trying to explain Y by X), it is sometimes a good idea to start with the right side of

$$(1) \quad f(x, y) = f_c(y | x) f_x(x)$$

(where $f(x, y)$ is the joint pdf of (X, Y) , $f_c(y | x)$ is the conditional pdf of $(Y | X = x)$, X being fixed to a value x , and $f_x(x)$ the marginal pdf of X .)

There are, however, also cases where the modelling is done directly from the left side of the factorization. Starting from the left side of (1), a linear relationship between X and Y is sometimes proposed directly as

$$(2) \quad Y = \alpha + \beta X + e$$

where the random error term, e , is uncorrelated with X , has expectation, $E(e) = 0$, and constant variance, $\text{var}(e) = \sigma_e^2$.

As clarified by Harald Cramér in his famous textbook from 1945, **Mathematical Methods of Statistics**, the relation (2) *should not be taken as an (additional) assumption* since it actually follows (is true) from any model for $f(x, y)$, (as long as moments up to second order exist) – which will be shown as follows:

(3) **Model:** Let (X, Y) have a joint distribution given by $f(x, y)$, such that expectations, variances and covariance exist. Write them as the following five population quantities

$$E(X) = \mu_X, \quad E(Y) = \mu_Y, \quad \text{var}(X) = \sigma_X^2, \quad \text{var}(Y) = \sigma_Y^2, \quad \text{and} \\ \text{cov}(X, Y) = \sigma_{XY}.$$

I will show that the statement (2) is automatically valid in model (3) without extra assumptions:

$$(4) \quad \text{Define}^1 \alpha \text{ and } \beta \text{ by: } \beta = \frac{\overset{\text{Def}}{\sigma_{XY}}}{\overset{\text{Def}}{\sigma_X^2}} \text{ and } \alpha = \mu_Y - \beta \mu_X$$

¹ Note that the definition of α and β , which are population quantities, corresponds exactly to the OLS estimates based on a sample of observations of (X, Y) .

(5) After this, define the rv, e , by: $e \stackrel{Def}{=} Y - \alpha - \beta X$

From this we have

$$(6) \quad Y = \alpha + \beta X + e$$

Now the rv, e , can be interpreted as an error term in the sense of (2), since we shall prove that $E(e) = 0$ and $\text{cov}(e, X) = 0$:

Proof of $E(e) = 0$ and $\text{cov}(e, X) = 0$:

We find from STAT1 –rules

$$E(e) = E(Y - \alpha - \beta X) = \mu_Y - \alpha - \beta \mu_X \stackrel{(4)}{=} \mu_Y - (\mu_Y - \beta \mu_X) - \beta \mu_X = 0$$

From this we get

$$\begin{aligned} \text{cov}(e, X) &= E(eX) - E(e) \cdot E(X) = E(eX) = E((Y - \alpha - \beta X) \cdot X) = \\ &= E(XY - \alpha X - \beta X^2) = E(XY) - \alpha \mu_X - \beta E(X^2) \end{aligned}$$

Substituting for α and β and using that $E(XY) = \text{cov}(X, Y) + E(X) \cdot E(Y) = \sigma_{XY} + \mu_X \mu_Y$, and $E(X^2) = \sigma_X^2 + \mu_X^2$, we get

$$\text{cov}(e, X) = \sigma_{XY} + \mu_X \mu_Y - \left(\mu_Y - \frac{\sigma_{XY}}{\sigma_X^2} \mu_X \right) \cdot \mu_X - \frac{\sigma_{XY}}{\sigma_X^2} (\sigma_X^2 + \mu_X^2) = \dots = 0 \quad \text{End of proof.}$$

So the conditions in (2) are automatically fulfilled.

We also need an expression for $\sigma_e^2 = \text{Var}(e)$. Having established that e and X are uncorrelated, we have

$$\sigma_Y^2 = \text{var}(Y) = \text{var}(\alpha + \beta X + e) = \beta^2 \text{var}(X) + \text{var}(e) = \beta^2 \sigma_X^2 + \sigma_e^2$$

Hence, introducing the correlation between X and Y , $\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$, we find

$$\sigma_e^2 = \sigma_Y^2 - \beta^2 \sigma_X^2 = \sigma_Y^2 - \frac{\sigma_{XY}^2}{\sigma_X^4} \sigma_X^2 = \sigma_Y^2 - \frac{\sigma_{XY}^2}{\sigma_X^2} = \sigma_Y^2 - \frac{\sigma_{XY}^2}{\sigma_X^2 \sigma_Y^2} \cdot \sigma_Y^2 = \sigma_Y^2 - \rho^2 \sigma_Y^2$$

or

$$(7) \quad \sigma_e^2 = \sigma_Y^2(1 - \rho^2)$$

(7) is a famous relationship that, among other things, offers an interpretation of ρ^2 (in the population) as a measure of the fraction of the total variance, σ_Y^2 , of Y that is explained by the linear relationship (2).

Note that another bonus that we get from (7) is a proof that $-1 \leq \rho \leq 1$ (assuming $\text{var}(Y) > 0$). This follows trivially since $\sigma_e^2 = \text{var}(e) \geq 0$ implies $\rho^2 \leq 1$.

2. Interpretation of (2)

Cramér derived (2) as the solution of a minimization problem. He showed that the values given of α and β , minimize the function, $Q(\alpha, \beta) = E[(Y - \alpha - \beta X)^2]$ (assuming (3) only), and that the minimizing values are unique². This represents one interpretation. I will give you another supplementary interpretation.

Write the linear part of (2) as $g(X)$, where g is the linear function, $g(x) = \alpha + \beta x$. As shown, this function is always well defined under model (3). A common misunderstanding appears to be a tendency to interpret $g(x)$ as the conditional expectation, $\mu(x) = E(Y | x)$, which is wrong in general. To claim that $g(x) = \mu(x)$, requires additional assumptions in model (3). In general $\mu(x)$ does not have to be linear at all, as examples 1 and 2 below, show. However, there is a certain relationship between the two functions $\mu(x)$ and $g(x)$. To see this, make first $\mu(x)$ random by replacing x by the rv X , leading to the rv, $\mu(X)$. Then, add and subtract this in Q , giving

$$\begin{aligned} Q(\alpha, \beta) &= E[(Y - g(X))^2] = E[(Y - \mu(X) + \mu(X) - g(X))^2] = \\ &= E[(Y - \mu(X))^2] + 2E[(Y - \mu(X))(\mu(X) - g(X))] + E[(\mu(X) - g(X))^2] \end{aligned}$$

Using the theorem of iterated expectations, we can show that the first term is equal to, $E[\sigma^2(X)]$, where $\sigma^2(x) = \text{var}(Y | x)$ is the conditional variance function.

[You may try to show this as an exercise. Write first $E[(Y - \mu(X))^2] = E[E((Y - \mu(X))^2 | X)]$ and use the two-step approach described in the lecture. On step one, find the function of x , $E((Y - \mu(X))^2 | X = x) = E((Y - \mu(x))^2 | x) \stackrel{\text{def}}{=} \text{var}(Y | x) = \sigma^2(x)$, noting that fixing X to the value x , turns, the rv $\mu(X)$ into the constant value, $\mu(x) = E(Y | x)$. On step two, replace x by the rv X and take expectations.]

Similarly, the second term in Q , can be shown to be 0. Hence, we can write

² Actually he showed (2) in a more general setting with an arbitrary number of explanatory variables, and not only one as here.

$$Q(\alpha, \beta) = E[\sigma^2(X)] + E[(\mu(X) - \alpha - \beta X)^2]$$

Since the minimization Q with respect to α and β does not affect the first term, we can now interpret the function $g(x)$ as the best linear approximation to $\mu(x)$ in an expected squared error sense.

3 Estimation of $g(x) = \alpha + \beta x$ in (2)

The joint distribution of (X, Y) , determined by the joint pdf, $f(x, y)$, may be called *the population distribution* (as suggested by Ragnar in the lecture), and we want information on this distribution based on a representative sample (data) from the population expressed as “ n observations³, $(x_1^o, y_1^o), (x_2^o, y_2^o), \dots, (x_n^o, y_n^o)$ of (X, Y) , sampled independently”. A model for this is the *iid* model where we consider the observations as observations of n random pairs, $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ assumed to be *iid* pairs. This means that we assume the n pairs to be independent random pairs, identically distributed as (X, Y) , and with the common pdf, $f(x, y)$.

Then, using the construction of (2) for each pair⁴, we get the (simple regression) model for data (known from several books),

$$Y_i = \alpha + \beta X_i + e_i, \quad i = 1, 2, \dots, n$$

where the error terms, e_1, e_2, \dots, e_n , are independent and identically distributed with expectation 0 and constant variance, $\text{var}(e_i) = \sigma_e^2$. This specification can be well treated (estimated) by the OLS method.

4. Example 1

We can construct the joint pdf, $f(x, y)$, for the random pair, (X, Y) using the factorization, $f(x, y) = f_c(y|x)f_x(x)$, (see (1)), where the two factors on the right can be modeled as we wish – independently of each other. For example, suppose that X is uniformly distributed over the interval $[0, 1]$ with pdf

$$f_x(x) = \begin{cases} 1 & \text{for } 0 \leq x \leq 1 \\ 0 & \text{for } x \text{ outside } [0, 1] \end{cases}$$

This implies that, $E(X) = \mu_x = 1/2$, and $\text{var}(X) = \sigma_x^2 = 1/12$.

³ The upper index o signifying that the observations are concrete numbers.

⁴ Note that, since all the pairs have the same distribution, α and β will be the same for each pair.

Fixing X to a number x between 0 and 1, we assume that $Y | (X = x)$ is normally distributed with expectation, $2x^2$, and variance 1 (in short $(Y | X = x) \sim N(2x^2, 1)$). This implies that the (true) regression function is $\mu(x) = E(Y | x) = 2x^2$ (well defined for x in $[0,1]$ only since the conditional pdf, $f_c(y | x)$, is not defined for x outside $[0,1]$).

[**Small technical point:** Note that, even if $f_c(y | x)$ is not defined for x outside $[0,1]$, we define (as is usual in practice) $f(x, y) = f_c(y | x)f_x(x) = 0$ for x outside $[0,1]$, since $f_x(x) = 0$ there. Hence, the joint pdf becomes

$$f(x, y) = \begin{cases} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y-2x^2)^2} & \text{for } 0 \leq x \leq 1 \text{ and } -\infty < y < \infty \\ 0 & \text{for any } (x, y) \text{ where } x \text{ is outside } [0,1] \end{cases}$$

Note that $f(x, y)$ is concentrated (i.e., > 0) in the strip which goes all along the y -axis determined by $0 \leq x \leq 1$ and 0 outside the strip. Note also that (X, Y) is *not* jointly normally distributed (the joint normal pdf is never 0 and would imply that both marginal distributions are normal).]

Hence the true regression function, $\mu(x)$, is not linear (but part of a parabola, see fig. 1 below). To determine the best linear approximation described in (2), we need as in (3), $E(X) = \mu_x = 1/2$, $E(Y) = \mu_y$, $\text{var}(X) = \sigma_x^2 = 1/12$, $\text{var}(Y) = \sigma_y^2$, and $\text{cov}(X, Y) = \sigma_{xy}$. The obvious tool here is the iterated expectation theorem:

$$\mu_y = E(Y) = E[E(Y | X)] = E[2X^2] = 2 \left[\text{var}(X) + (E(X))^2 \right] = 2 \left(\frac{1}{12} + \frac{1}{4} \right) = \frac{2}{3}$$

$$\sigma_{xy} = \text{cov}(X, Y) = E(XY) - \mu_x \mu_y = E(XY) - \frac{1}{3}$$

$$E(XY) = E[E(XY | X)] = E[X \cdot E(Y | X)] = E[X \cdot 2X^2] = 2E[X^3]$$

[**Note.** To understand this manipulation, it is best to use the two-stage approach described in the lecture. **Step 1:** Find the conditional function behind the inner expectation first. $c(x) = E(XY | X = x) = E(xY | X = x) = x \cdot E(Y | x)$. This works since the value x is just a constant in the distribution of $Y | (X = x)$. Hence, $c(x) = xE(Y | x) = 2x^3$. **Step 2:** Replace x by the rv X in $c(x)$ and take expectation: $c(X) = X \cdot E(Y | X) = 2X^3$. The theorem of iterated expectations tells us that taking the expectation of $c(X)$ gives us $E(XY)$.]

$$E(XY) = 2E(X^3) = 2 \int_0^1 x^3 f_x(x) dx = 2 \int_0^1 x^3 dx = 2 \cdot \frac{1}{4} = \frac{1}{2}$$

Hence

$$\sigma_{xy} = E(XY) - \frac{1}{3} = \frac{1}{2} - \frac{1}{3} = \frac{1}{6}$$

$$\sigma_Y^2 = \text{var}(Y) = E[\text{var}(Y | X)] + \text{var}[E(Y | X)] = E[1] + \text{var}[2X^2] = 1 + 4 \cdot \text{var}(X^2)$$

Hence

$$\sigma_Y^2 = 1 + 4 \left(E[X^4] - (E[X^2])^2 \right) = 1 + 4 \left(\int_0^1 x^4 \cdot 1 dx - \left(\int_0^1 x^2 \cdot 1 dx \right)^2 \right) = 1 + 4 \cdot \frac{4}{45} = 1 + \frac{16}{45}$$

We can now find the best linear approximation, $g(x) = \alpha + \beta x$, to $\mu(x) = E(Y | x)$, where

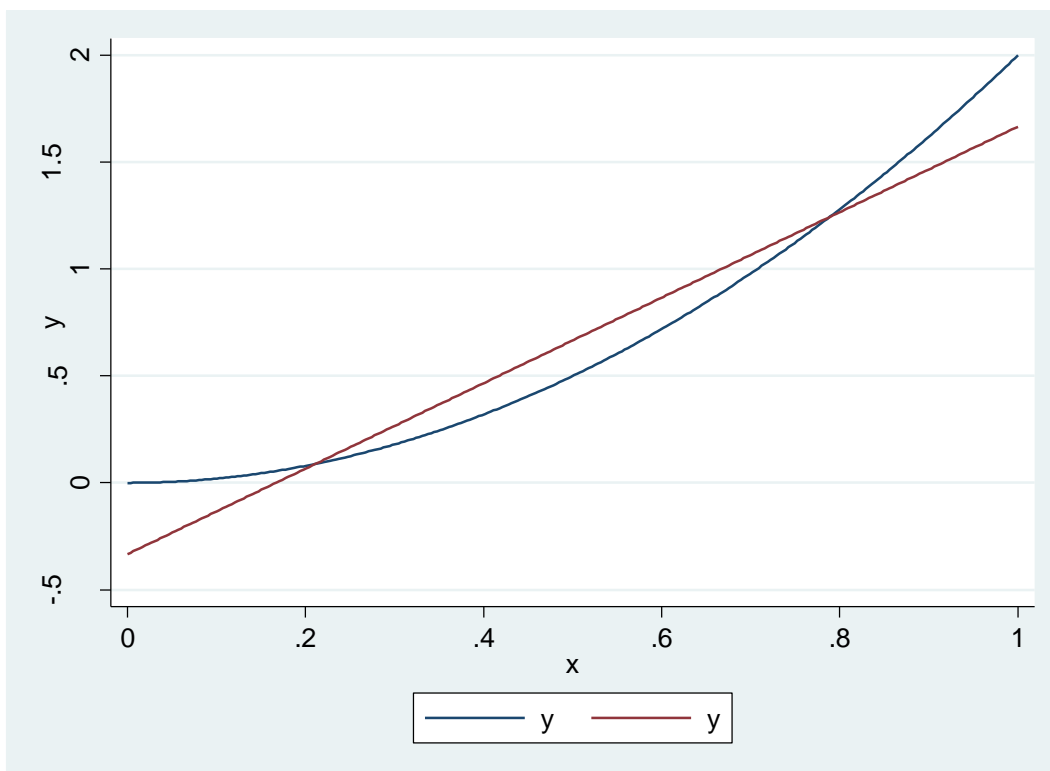
$$\beta = \frac{\sigma_{XY}}{\sigma_X^2} = \frac{1/6}{1/12} = 2$$

$$\alpha = \mu_Y - \beta \mu_X = \frac{2}{3} - 2 \cdot \frac{1}{2} = -\frac{1}{3}$$

or

$$g(x) = -\frac{1}{3} + 2x$$

Figure 1 $g(x) = \alpha + \beta x$ and regression function $\mu(x) = E(Y | x)$



Stata command: twoway (function y=2*x^2, range(0 1)) (function y=-1/3+2*x, range(0 1))

Example 2 Exercise

A. Repeat the calculations as in example 1, now assuming

(i) $X \sim$ uniformly distributed over $[0, 2]$

[implying, $f_X(x) = 1/2$ for $0 \leq x \leq 2$, $\mu_X = E(X) = 1$, and $\sigma_X^2 = \text{var}(X) = 1/3$]

(ii) $Y | x \sim N(2(x-1)^2, x^2)$

[implying, $\mu(x) = E(Y | x) = 2(x-1)^2$ and $\sigma^2(x) = \text{var}(Y | x) = x^2$]

Hint. Verify that $\mu_Y = \frac{2}{3}$, $\sigma_Y^2 = \frac{76}{45}$, $\sigma_{XY} = -\frac{1}{3}$.

Developing σ_Y^2 , you may need to find $\text{var}[(X-1)^2] = E[(X-1)^4] - (E[(X-1)^2])^2$.

Note that, for example, $E[(X-1)^4] = \int_0^2 (x-1)^4 f_X(x) dx = \int_0^2 (x-1)^4 \cdot \frac{1}{2} dx = \frac{1}{5}$

B. Find the best linear approximation, $g(x) = \alpha + \beta x$, to the true regression, $\mu(x) = E(Y | x)$, and plot both functions in the same graph.

C. How would you go about to simulate (make the computer draw) n independent observations of the random pair, (X, Y) ? (**Hint:** Utilize the right side of (1).)

(Try it (!), using (e.g.) Excel or Stata for $n = 50$, and make a scatter plot.)