

When to use regression analysis (linear model)?

回歸

What data? What objective?

meaning?
LNp. 2-4

Regression: a statistical tool for investigating the "linearity relationship" between x and y .

因果關係

could have many x variables

效應

causal relationship: examine the effects of x on y , i.e. how the changes in x result in the change in y (note that the statement implicitly assumes causal relationship), often seen in

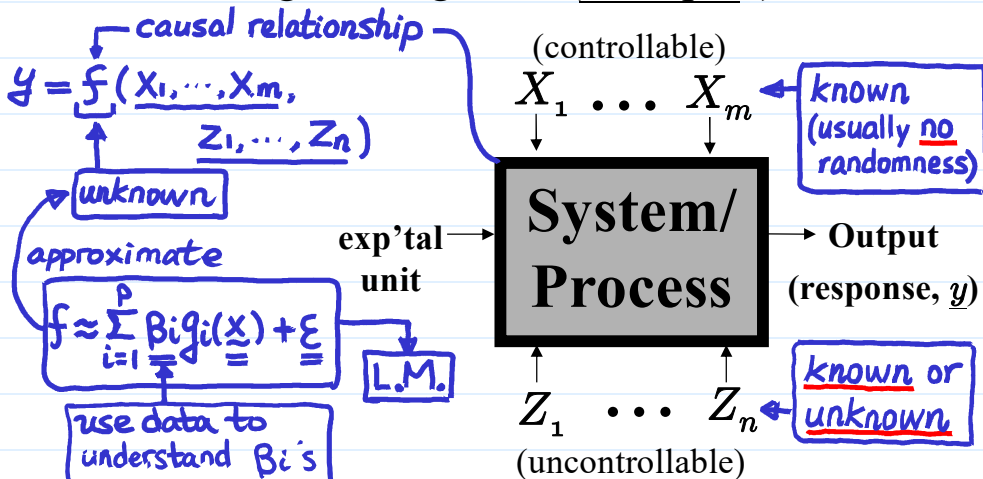
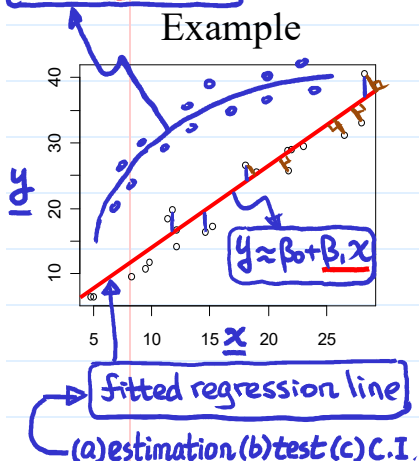
x : 因
 y : 果

not easy to build purely based on data analysis

- DOE data (example?)
- physical/chemical/engineering data (example?)

functional relationship

$$y \approx \beta_0 + \beta_1 x + \beta_2 x^2$$



Even where no sensible causal relationship exists between x and y , we may wish to relate them by some sort of mathematical equation (rationale: sample from a multivariate normal population), often seen in

association
 x & y
聯動

X usually not random in DOE

strong association
~~causal relation~~

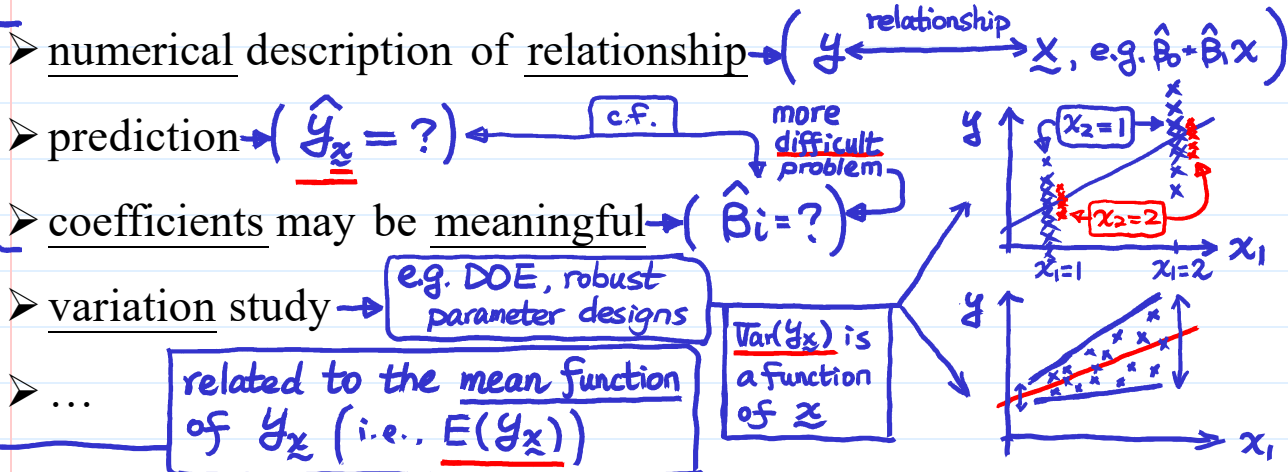
- social data (example?)
- economical data (example?)
- medical data (example?)

X & Y : random variables
 $(X, Y) \sim$ joint distribution
 $Y|X \rightarrow$ linear model
(future lecture: sampling model)

- shoe size, reading
- smoking, cancer

In some cases, what variable should be treated as the response & what should be the predictors could be a difficult problem.

Q: Why develop the fitted line (fitted model)? Why it is useful? LNp. 2-8



• data type in regression analysis and some terminologies

Recall continuous vs. discrete random variables

A **反應變數** $\{response, output, dependent\}$ variable Y is modeled or explained by p effects/functions of m **解釋變數** $\{predictor, input, independent, regressor\}$ variables X_1, \dots, X_m

Note: All "observed" data are discrete

- Y : "approximately" continuous **c.f.** $\left[\begin{matrix} 氣候 \\ 溫度 \\ 人體溫 \\ 人數 \\ 學校人數 \\ 中獎人數 \end{matrix} \right]$ **c.f.** $\left[\begin{matrix} normal \\ fever \end{matrix} \right]$
- X_1, \dots, X_m : continuous and discrete (quantitative), categorical (qualitative) (**Q**: example?)

of effects

- $p=1$, simple regression; $p>1$, multiple regression
- X_1, \dots, X_m
 - all quantitative \Rightarrow multiple regression
 - quantitative+qualitative \Rightarrow analysis of covariance
 - all qualitative \Rightarrow analysis of variance (ANOVA)

All \in Linear models

$e.g. \rightarrow y = \beta_0 + \beta_1 x + \epsilon$

more than one Y , multivariate regression \rightarrow multivariate data analysis

• linear model:

$$E(Y|X_1, \dots, X_m) = \sum_{i=0}^p \beta_i \cdot g_i$$

$$Var(Y|X_1, \dots, X_m) = Var(\epsilon) = \sigma^2$$

$$Y = \sum_{i=0}^p \beta_i \cdot g_i(X_1, \dots, X_m) + \epsilon$$

規律 (signal) \leftarrow deterministic component \rightarrow mean function \leftarrow **隨機** (noise) \leftarrow error: random component \rightarrow variance function

$E(\epsilon) = 0$
 $Var(\epsilon) = \sigma^2$

X_1, \dots, X_m are regarded as deterministic, i.e., no random phenomenon (when they are random variables, regard the linear model as conditional on X_1, \dots, X_m)

how to determine g_i 's?

- $g_0(X_1, \dots, X_m), \dots, g_p(X_1, \dots, X_m)$: known functions of X_1, \dots, X_m , called effects
- (unknown) parameters β_0, \dots, β_p enter linearly
- variation due to random error only appears on y -axis

know the form of effects (g_i 's) but don't know their magnitudes (β_i 's)

• Rationale: a general model for the relationship between Y and X_1, \dots, X_m is:

$$Y = f(X_1, \dots, X_m) + \epsilon, \text{ where } f \text{ is unknown and arbitrary}$$

Why? consider the Taylor expansion of $f \rightarrow f = \sum_{i=0}^{\infty} \beta_i x^i$

Note: # of parameters in f is infinite, usually do not have enough data to estimate f directly (globally), we have to assume that it has some more restricted form

local approximation of f may be achievable by a linear model **with finite # of parameters**

Note: Because the predictors can be transformed and combined in any way, linear models are actually very flexible. **e.g. g_i 's: $x_1^t, x_1 x_2, e^{x_1+x_2}, \log(\frac{x_1}{x_2}), \dots$**

• Example:

Data

	Y	X ₁	X ₂	X ₃
car1	•	•	•	•
car2	•	•	•	•

causality or association?

Y: fuel consumption of a particular model of car

X₁: weight of the car (continuous? discrete? categorical?),

X₂: horse power (continuous? discrete? categorical?),

X₃: number of cylinders (continuous? discrete? categorical?)

X₄: 廠牌, Honda (1), Ford (2), Toyota (3), ... (categorical)

➤ a general model: $Y = f(X_1, X_2, X_3) + \epsilon$, where f is unknown & arbitrary

➤ possible statistical models ("local approximation" of f):

many

parameter

$Y = \beta_0 \underline{1} + \beta_1 \underline{X_1} + \beta_2 \underline{X_2} + \beta_3 \underline{X_3} + \epsilon$ ⇒ a linear model

$Y = \beta_0 + \beta_1 \underline{X_1^2} + \beta_2 \underline{\log(X_2)} + \beta_3 \underline{X_1 X_3} + \epsilon$ ⇒ a linear model

$Y = \beta_0 + \beta_1 \underline{X_1^{\beta_2}} + \epsilon$ ⇒ not a linear model

$Y = \beta_1 X_1^{\beta_2} \epsilon$ ⇒ not a linear model, but can be a linear model after taking log

$\log(Y) = \log(\beta_1) + \beta_2 \log(X_1) + \log(\epsilon)$

Y' β_i' β ε'

❖ Reading: Faraway (2005, 1st ed.), 1.3, 2.1

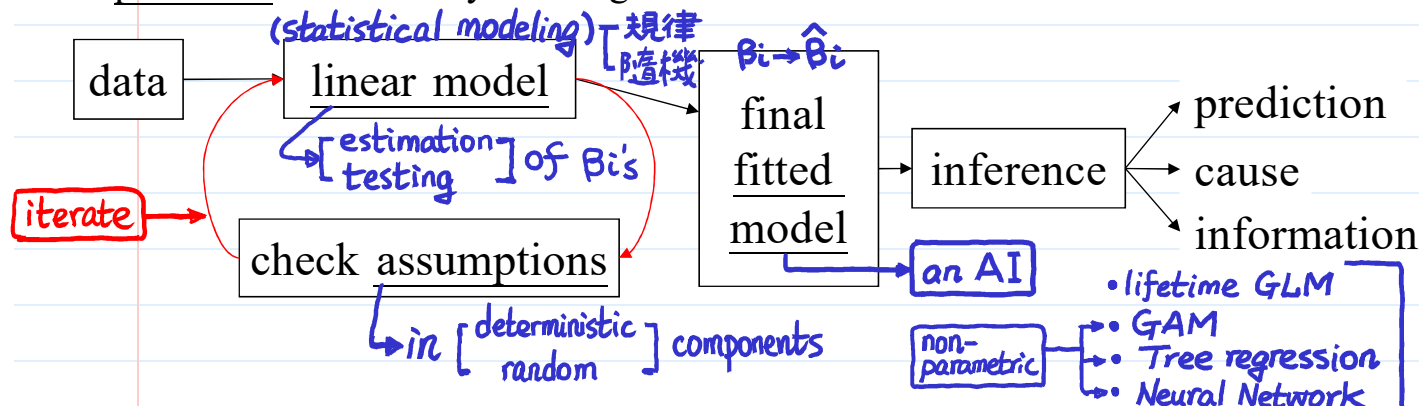
❖ further reading: D&S, 1.2 (meaning of linear model)

some possible objectives of regression analysis

LNp.2-2

1. prediction of future observations
2. assessment of the effect of, or relationship between, explanatory variables on the response
 ↑ e.g. DOE in LNp.2-1
3. a general description of data structure
 ← e.g. association between some variables (LNp.2-2)
4. ...

• procedure of data analysis in regression



• extensions exist to handle (1) multivariate responses, (2) binary responses (logistic regression), (3) count responses (Poisson regression), (4) ...

- (1) more than one Y (multivariate analysis) (2) $Y_x \sim \text{binomial}(n_x, p_x)$ (3) $Y_x \sim \text{Poisson}(\lambda_x)$
- GLM

History of regression

- ① linear structure
- ② data with random errors
- ③ parameter estimation

Tobias Mayer (1750), made numerous observations to the libration of the moon with the purpose of determining the characteristics of the moon's orbit, and proposed the method of averages

$$Y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1}$$

$$Y_{p \times 1}^* = X_{p \times p}^* \beta_{p \times 1}$$

$$\Rightarrow \beta = (X^*)^{-1} Y^*$$

Adrien Marie Legendre (1805), developed the method of least square generalized → Gauss-Markov Thm (future lecture) ← future lecture

Carl Friedrich Gauss (1809), claimed to have developed the method a few years earlier, and showed later least square is the best when the errors are normally distributed. *no need to get too pessimistic about mediocrity*

Sir Francis Galton (1875), coined the term "regression to mediocrity"

Family heights data

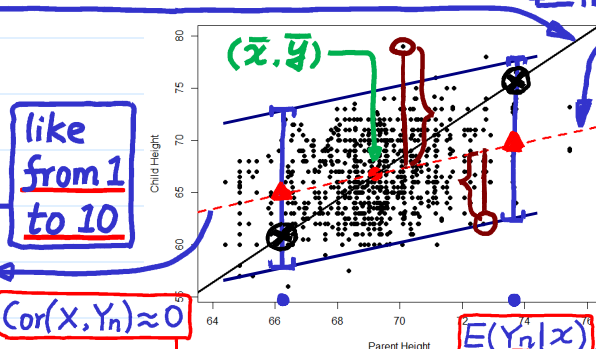
taller (shorter) parent
 ⇒ taller (shorter) child, but not as tall (short) as parent.

$$\left(\frac{y - \bar{y}}{\hat{\sigma}_y}\right) \approx 1 \left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right)$$

Standard score

$$\left(\frac{y - \bar{y}}{\hat{\sigma}_y}\right) = \hat{\rho} \left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right)$$

$\hat{\rho} < 1$



least square line
 need to consider the influence of ϵ

zero $\left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right) \rightarrow$ 1st $\hat{\rho} \left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right) \rightarrow$ 2nd $\hat{\rho}^2 \left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right) \rightarrow \dots \rightarrow$ nth $\hat{\rho}^n \left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right)$

$\hat{\rho}^n \approx 0$, when n large $\Rightarrow \left(\frac{y_n - \bar{y}_n}{\hat{\sigma}_{y_n}}\right) \approx 0 \Rightarrow y_n \approx \bar{y}_n$

So, $\hat{\sigma}_{y_n} \approx 0$? No.

← This so called regression effect is not a natural law. It is actually a mathematical result.

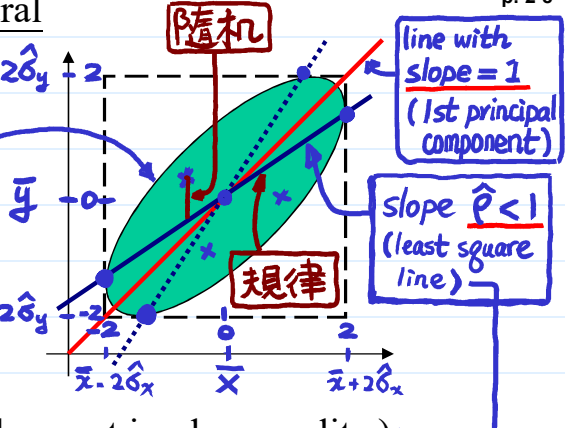
In many applications, regression effect is not of interest. (e.g., different types of variables)

$$\left(\frac{y - \bar{y}}{\hat{\sigma}_y}\right) = \hat{\rho} \left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right)$$

$-1 < \hat{\rho} < 1$

$$\left(\frac{x - \bar{x}}{\hat{\sigma}_x}\right) = \hat{\rho} \left(\frac{y - \bar{y}}{\hat{\sigma}_y}\right)$$

determined by the covariance matrix of (X, Y), LNp.4-5



➤ Misinterpretation (Note. Regression effect does not imply causality.)

check LNp.2-2

- IQ of a couple (husband, wife)
 - h: response, w: predictor
 - w: response, h: predictor
 - Ineffective treatment, measurements before and after treatment
 - Punishment works and reward does not.

should use 2-sample comparison or paired comparison

➤ Warm and Cold → pro/con

- Family height data
- Josef Stalin: "A single death is a tragedy; a million deaths is a statistic!" (一個人的死是一場悲劇，百萬個人的死只是一個統計數據！)

❖ Reading: Faraway (2005, 1st ed.), 1.4

❖ Further reading: D&S, 1.8;

Cathy O'Neil (2016), Weapons of Math Destruction (中譯: 大數據的傲慢與偏見).