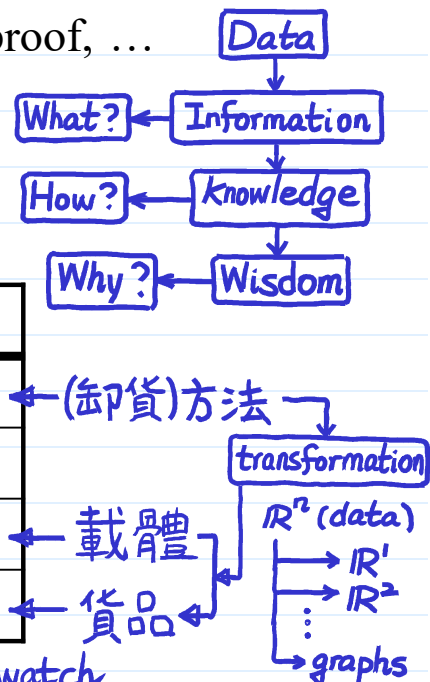## Question

# What is Statistics?

- A branch of <u>math</u> --- calculation, derivative, proof, …

- A collection of many <u>statistics</u> (formula)

- A useful <u>tools</u> for <u>extracting</u> <u>information/knowledge</u> from the <u>data</u>

| 哈利波特 | Real Life |
|---|---|
| 占卜學 | Statistics |
| 崔老妮 | Statisticians |
| <u>水晶球</u> | <u>Data</u> |
| 未來的資訊 | Information |

*aim of statistics*: provide *insight* by means of *data*

## **Basic Procedures of Statistics**

- Statistics divides the study of data into *five* steps:



- **Q**: What is a <u>statistical model</u>?
- $X_1, …, X_n$ (<u>random variables</u>) $\sim$ joint cdf $F_X$/pdf $f_X$/pmf $p_X$ with parameters $\Theta$

Data: $X_1, …, X_n$

Problem Formulation & <u>Modeling</u> (conceptual )

Data Collection

Statistical <u>Modeling</u> (empirical)

Data Analysis

Decision Making

Transformations
$g_1(X_1, …, X_n),$
…,
$g_k(X_1, …, X_n)$

$X_1, …, X_n \longrightarrow \Theta$

Extract Information

# 1. Problem formulation & modeling (conceptual approach)

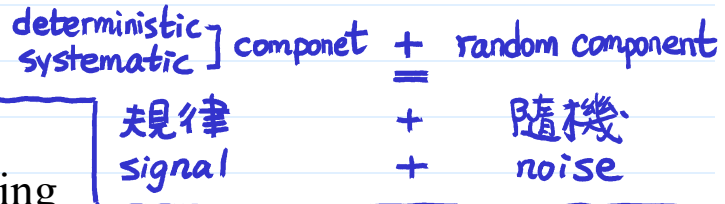*— key*

*problem*
*↓*
*formulation*
*↓*
*statistical problem*

➢ Problem <u>formulation</u>: use statistical/probabilistic/ mathematical language to "clearly" define the problem and the objective of study

*domain knowledge*

*key: good surrogate*

➢ modeling (<u>conceptual</u> approach): use the <u>information</u> that we possessed *prior to obtaining data* to develop a *representation of the underlying system*, also account for <u>uncertainty</u> in data

*deterministic systematic ] componet + random component*
*規律      +   隨機*
*signal    +   noise*

## 2. Data collection: producing *representative* data for drawing correct information

*有代表性的*

➢ survey sampling
   (抽樣調查)

➢ design of experiment
   (實驗設計)

➢ observational data

*Data*
① *may contain wrong information,*
② *may not contain useful information about the problem.*

## 3. Statistical <u>modeling</u> (empirical approach): use <u>empirical information</u> contained in the <u>data</u> to build a model or to <u>justify/adjust</u> the (conceptual) model developed in **1.**, also account for uncertainty in data

➢ a statistical model is a description of the <u>joint distribution</u> of data

*# of parameters ≠ ∞*

*deterministic component (規律)*

➢ a statistical model may contain the following components:
   • <u>nonparametric</u> component
   • <u>parametric</u> component: (<u>fixed</u>, <u>random</u>) effects
   • <u>distribution</u> component

*random component (隨机)*

*c.f.*

*Semi-parametric*

## 4. data analysis: mining <u>information</u> from data

*objective probability*

➢ <u>graphical</u> methods

*about parameters (規律)*

*subjective probability*

➢ <u>numerical</u> methods
   • (<u>point</u>, <u>interval</u>) estimation
   • <u>hypothesis testing</u>

## 5. Inference/decision making: drawing <u>conclusions</u> & <u>answering</u> questions based on <u>results</u> obtained in **4**.

• Example (from Gilchrist, *Statistical Modelling*, 1984):

> "A range of problems related to the <u>positioning</u> of <u>stores</u> and the planning of <u>delivery routes</u> requires information on the <u>distances by road</u>, *y*, <u>between different places</u>. Where a <u>large number</u> of such places are involved, finding these distances by driving or by direct measurement along the roads on a map is <u>time-consuming</u>."

*[handwritten: Q: problem "clearly" defined?]*

*[handwritten: Problem: how to measure road distance y btw any 2 stores.]*

*[handwritten: formulation]*

*[handwritten: a statistical problem]*

> "To avoid this problem, the usual approach is to relate the <u>road distances</u> *y* to the <u>straight line</u> distance, denoted by *x*, as measured using a scale map. This <u>relationship</u> will be expressed <u>mathematically</u> and will enable us to <u>predict</u> a value of *y* given a corresponding <u>value of x</u>. This relationship will be our <u>quantitative model</u> of the situation. The fundamental question is: *how do we obtain this relationship (model)*."

*[handwritten: $\hat{y} = f(x)$]*

*[handwritten: unknown]*
*[handwritten: • how can we understand it using data?]*
*[handwritten: • do we have some prior information about f?]*

Let's assume the following conditions (are they reasonable?):

a)  $x=0 \Rightarrow y=0$

*[handwritten: x=12, y=12]*  *[handwritten: x=12, y>12]*

b)  If there is a <u>straight road</u> between two points, then $x=y$; otherwise, $y \geq x$

c)  Generally, *y* should <u>increase</u> with *x*. However, because of <u>randomness</u> in road patterns, places with <u>same</u> *x*'s may have <u>different</u> *y*'s.

d)  Under similar situations, e.g. urban roads, the <u>form</u> of the relationship should not depend strongly on the <u>distances</u> involved, i.e., if *x* is, say, <u>doubled</u>, we would expect *y* is also <u>approximately</u> doubled.

Consider the following relationships (models):

*1.*  $y=x$                              [satisfies a) and d), but not b) or c)]

*2.*  $y=x+\varepsilon$, $\varepsilon$: random component       [now allows c), but not b)]

*3.*  $y=\alpha+x+\varepsilon$, $\alpha$: a constant          [helps with b), but a) fails]

*[handwritten: systematic component]*  *[handwritten: random component]*

*4.*  $y=\beta x+\varepsilon$, $\beta$: a constant $\geq 1$         [satisfies all four conditions. <u>true</u>?]

*5.*  <u>distribution assumption</u> can be added on the $\varepsilon$ in *4*, e.g., <u>$\varepsilon \sim N(0, \sigma^2)$</u>

**Note**: The above (<u>conceptual</u>) model is derived <u>without</u> any <u>data</u> provided.

<u>Problem formulation</u>: <u>Estimate</u> and <u>test</u> <u>parameters</u> in $y=\beta x+\varepsilon$, where $\beta \geq 1$

*[handwritten: a statistical problem "clearly" defined.]*  *[handwritten: variables in data set]*

## Some Notes in Problem formulation & modeling (conceptual approach)

— statistics: 輔助科學

⊙ understand the physical/social/political/biological/medical/... background to avoid the missing of important conditions that should be included in model

domain Knowledge

e.g. Cox proportional hazard model

- understand the objective

semi-parametric ─ ┌ parametric part (規律 of interest)
                  └ nonparametric part (規律 not of direct interest)

- make sure you know what the client wants

- state the problem in "statistical language"

> **Albert Einstein**. *The formulation of a problem is often more essential than its solution which may be merely a matter of mathematical or experimental skill.*

---

p. 1-8

Example (cont.):

- the collected data are given in the tabular. Is it a "representative" data set?

  e.g. record data on some day

- observational or experimental data?

  e.g., uniformly choose stores on the map.

large leverage

| y | x |
|------|------|
| 10.7 | 9.5 |
| 6.5 | 5 |
| 29.4 | 23 |
| 17.2 | 15.2 |
| 18.4 | 11.4 |
| 19.7 | 11.8 |
| 16.6 | 12.1 |
| 29 | 22 |
| 40.5 | 28.2 |
| 14.2 | 12.1 |
| 11.7 | 9.8 |
| 25.6 | 19 |
| 16.3 | 14.6 |
| 9.5 | 8.3 |
| 28.8 | 21.6 |
| 31.2 | 26.5 |
| 6.5 | 4.8 |
| 25.7 | 21.7 |
| 26.5 | 18 |
| 33.1 | 28 |

- **Q**: If you can design the experiment, what are the data collection issues that should be concerned in the example?
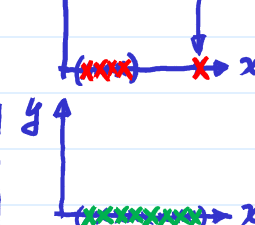
- Consider the following situations:

  ➢ if there are hundred/thousand of places, how to choose a small number of appropriate locations? geometrically uniform allocation? stratified sampling?

  ➢ what if there are many routes that link any two places? replication required?    ➤ can be used to understand "information" about the variation caused by routes

  ➢ who should be assigned to measure these $y$'s by driving? randomization? blocking?

## Some Notes in Data Collection

- are the data observational or experimental?

- how to collect a representative data?

- is there non-response?　　*It's also informative.*

- are there missing values?　　*missing information.*
  *MCAR, MAR, MNAR*

- qualitative or quantitative?
  類別型　　　　連續型

- how are the data coded?

- what are the units of measurement?

- beware of data entry errors　← *Data sanity check*
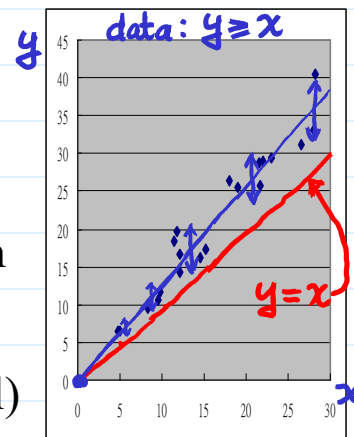  *Data cleaning*

Example (cont.):

- What empirical model will you suggest after examining the plot?

- should empirical model be identical to conceptual model?

- if the plot (or numeric analysis) reveals different patterns ...

  ➢ what if you find curvature or jump relationship existing between $x$ and $y$?　　$y = \beta x + \varepsilon, \beta \geq 1$

  ➢ what if you find non-constant variance?

  how should the conceptual model be adjusted?
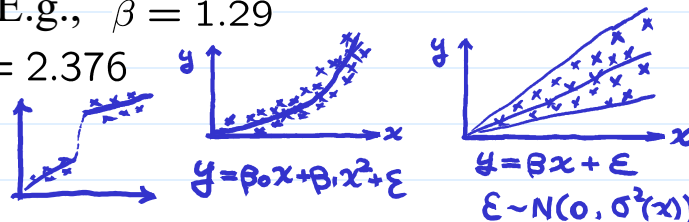
- graphic analyses offer vivid and intuitive perception

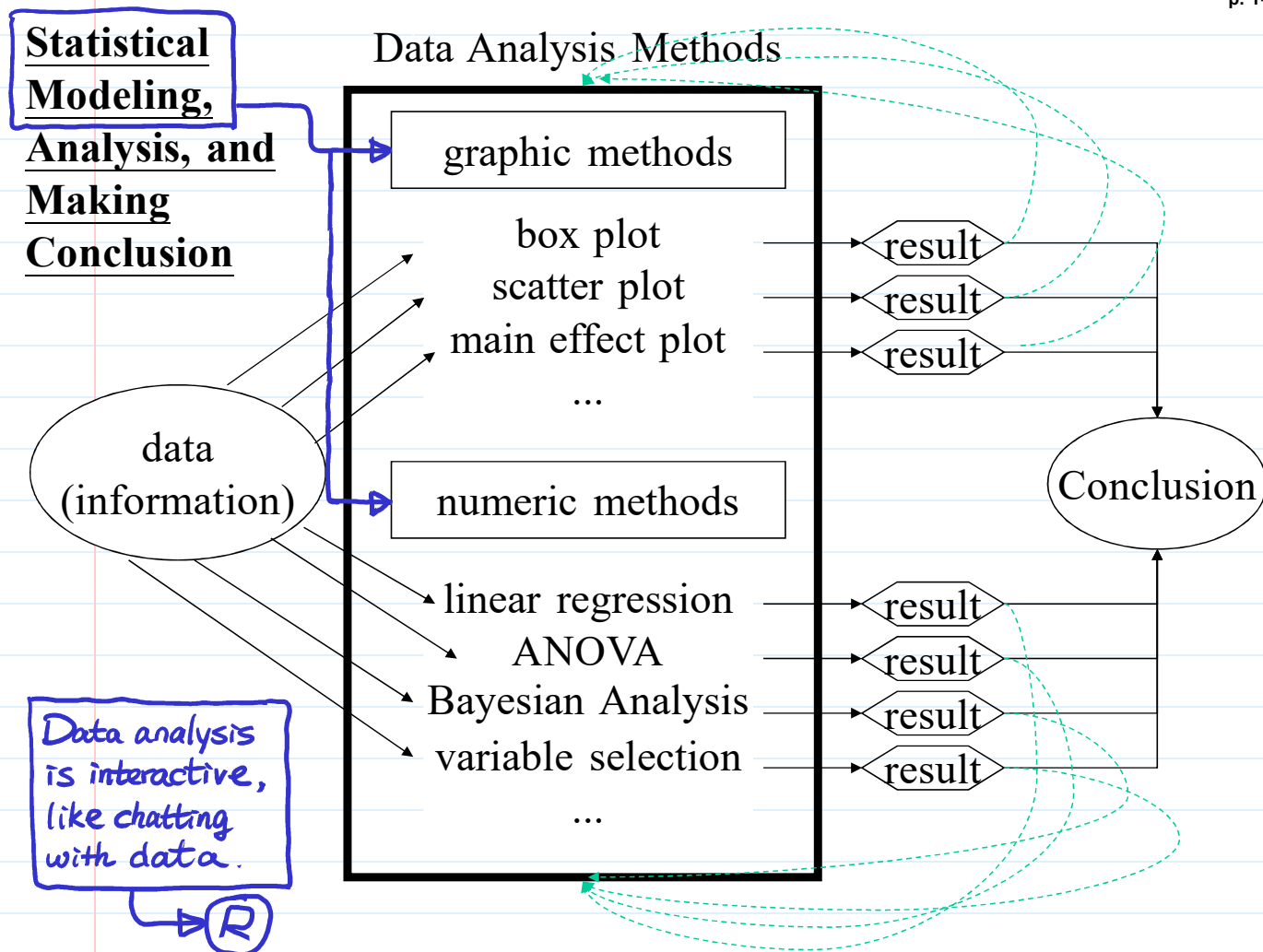- numeric analyses present numeric summaries (such as estimation and testing of parameters in the model) for making concrete conclusions. E.g., $\hat{\beta} = 1.29$ and is significant in $t$-test, and $\hat{\sigma} = 2.376$

*final fitted model*

⊙ Conclusion: $\hat{y} = 1.29x$
(or offer confident interval of $\hat{y}$)

*data: y ≥ x*
*y = x*

$y = \beta_0 x + \beta_1 x^2 + \varepsilon$

$y = \beta x + \varepsilon$
$\varepsilon \sim N(0, \sigma^2(x))$

**Statistical Modeling, Analysis, and Making Conclusion**

Data Analysis Methods

graphic methods

box plot
scatter plot
main effect plot
...

result
result
result

numeric methods

linear regression
ANOVA
Bayesian Analysis
variable selection
...

result
result
result
result

data (information)

Conclusion

*Data analysis is interactive, like chatting with data.*  ⓡ
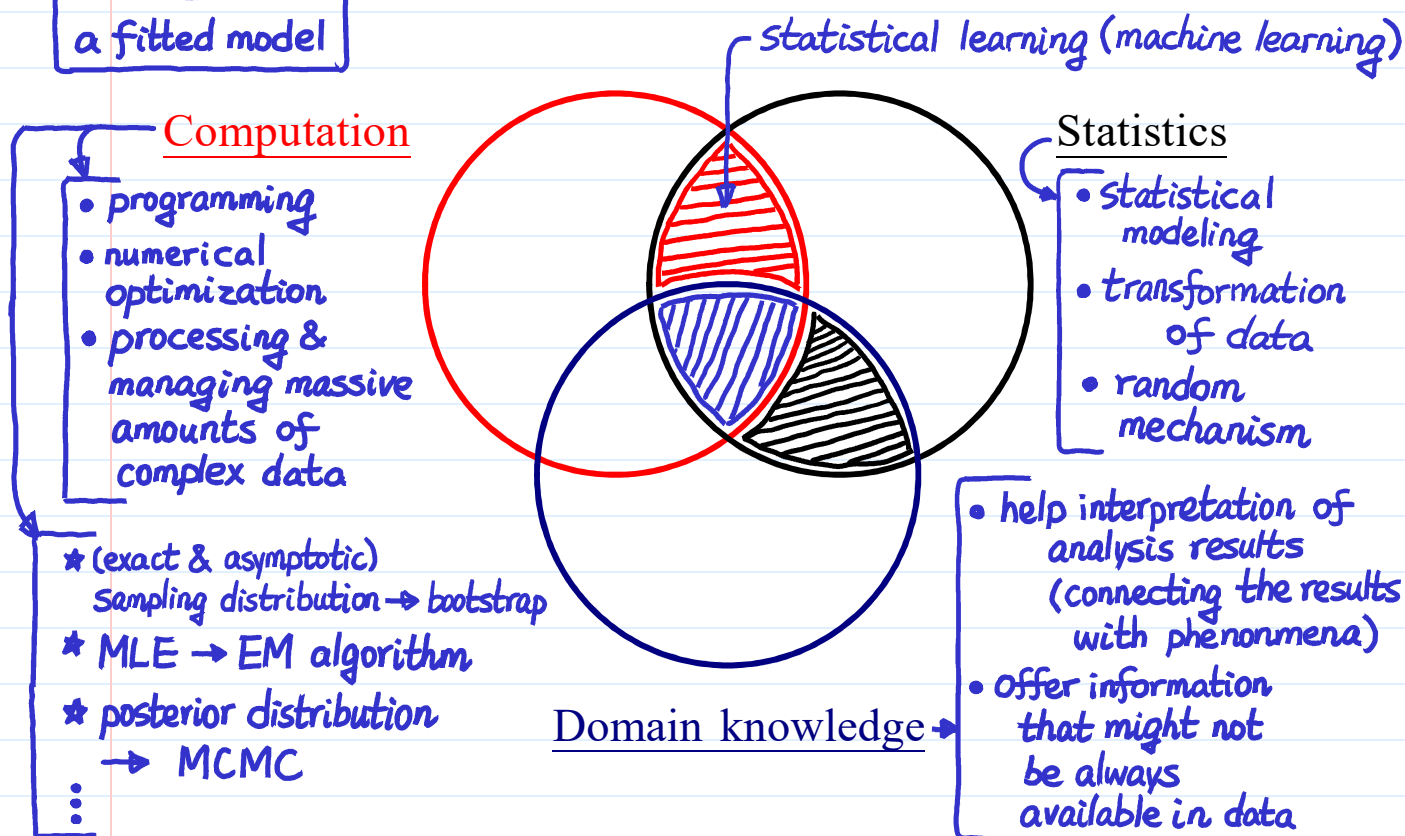
---

**Some Notes in Statistical modeling, Data analysis, and Decision making**

- If possible, most <u>available analysis methods</u> should be performed.   *Data analysis ⟺ projection from different angles*

- <u>Assumptions</u> and analysis results between <u>different</u> <u>methods</u> could be (slightly) <u>different</u>
  *→ hidden in statistical models*

- Data analysis is inherently <u>interactive</u>

- <u>Conclusions</u> should be summarized based on <u>consistent</u> results.  *→ level of evidence*

- <u>Important information</u> usually <u>consistently appear</u> in the results of every methods

- <u>quantitative</u> (定量) and <u>qualitative</u> (定性) conclusions

# A successful data analysis usually requires
# a mix of the three components:

an AI is
a fitted model

*statistical learning (machine learning)*

Computation

Statistics

- programming
- numerical optimization
- processing & managing massive amounts of complex data

- statistical modeling
- transformation of data
- random mechanism

★ (exact & asymptotic) sampling distribution → bootstrap

★ MLE → EM algorithm

★ posterior distribution → MCMC

⋮

Domain knowledge →

- help interpretation of analysis results (connecting the results with phenonmena)
- offer information that might not be always available in data

❖ **Reading**: Faraway (2005, 1st edition), 1.1

required

❖ **Further reading**: *D&S : Draper and Smith (1998)*
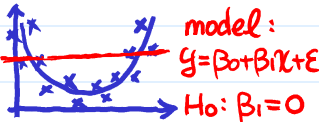
optional

- Statistical modelling (Gilchrist, 1984)

- Statistics: a guide to the unknown (edited by Tanur et al., 1972, 1978, 1989; Peck et al., 2005)

- Applied statistics: principles and examples (Cox & Snell, 1981)

❖ **Some other reading:**

- Lewis (2004), Moneyball (中譯：魔球).

- Kahneman (2011), Thinking, Fast and Slow (中譯：快思慢想).

- Silver (2012), The Signal and the Noise (中譯：精準預測).

## What aspects you should focus on in this course?

1. **Understand** analysis methods

   $H_0 \cup H_1$: collection of all the models (parameters) considered in the test.

   - objective is ...?

     model:
     $y = \beta_0 + \beta_1 x + \varepsilon$
     $H_0: \beta_1 = 0$

   - for an estimator (parameter), what's its meaning?

   - for a test, what are its $H_0$ and $H_1$?   collection of models

   - how to find statistically significant results in outputs?

   - assumptions and limitations in a statistical model?

   - ...

     not only p-value > or < 0.05
     level of evidence (LNp. 1-12)

2. **Interpretation**: for those significant results, how to interpret them in the language that your clients use

3. **How to implement** the analysis method in softwares, such as R, Splus, SAS, ...?