# Introduction to Statistics

Berlin Chen
Department of Computer Science & Information Engineering
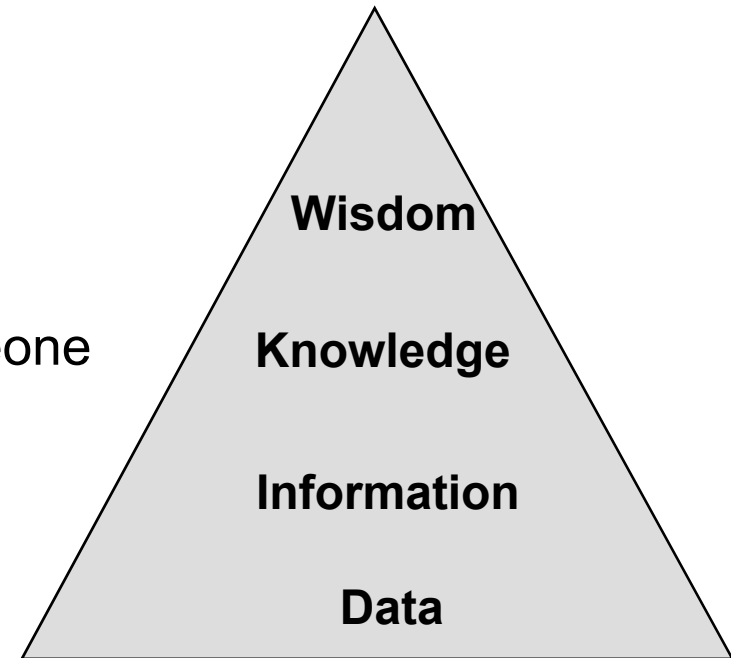National Taiwan Normal University

# What is Statistics?

- Statistics is the field of study concerned with the collection, analysis, and interpretation (making decisions on) of uncertain data
  - E.g., the explanation of social or economic trends through the analysis of data

- Or, in more common usage, statistics refers to numerical facts of the data
  - E.g., the age of a student, the allowance of a student, the height of a student, etc.

- Another definition: Statistics is the science of conducting studies to collect, organize, summarize, analyze, and draw conclusions from data
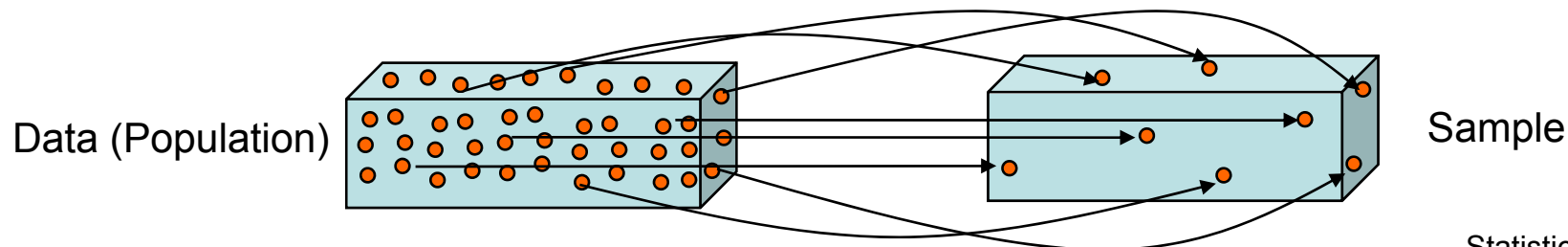
統計學：“以偏概全”＋“有所本”??

# Information Hierarchy

- ## Data
  - The raw material of information

- ## Information
  - Data organized and presented by someone

- ## Knowledge
  - Information read, heard or seen and understood

- ## Wisdom
  - Distilled and integrated knowledge and understanding

**Wisdom**

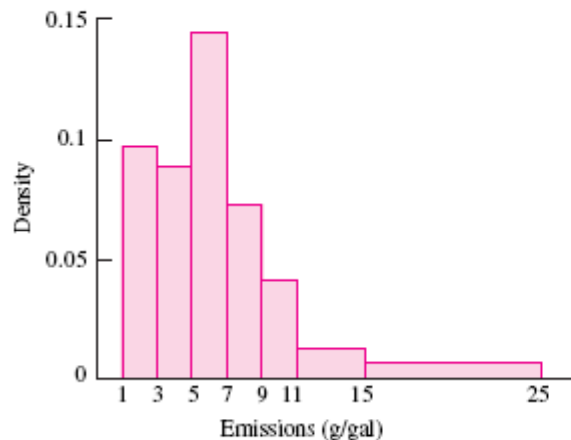**Knowledge**

**Information**

**Data**

# Types of Statistics (1/4)

- Broadly speaking, statistics can be divided into two areas
  - Descriptive statistics (敘述統計學)
  - Inferential statistics (推論統計學)

- Descriptive Statistics
  - To be concerned with the methods of collecting data and of summarizing clearly the basic information they contain
    - Collecting data refers to sampling, i.e., choosing a subset of data (a sample)
    - Summarizing data refers to organizing, displaying, and describing data by tables, graphs, and summary measures

Data (Population)
Sample

# Types of Statistics (2/4)



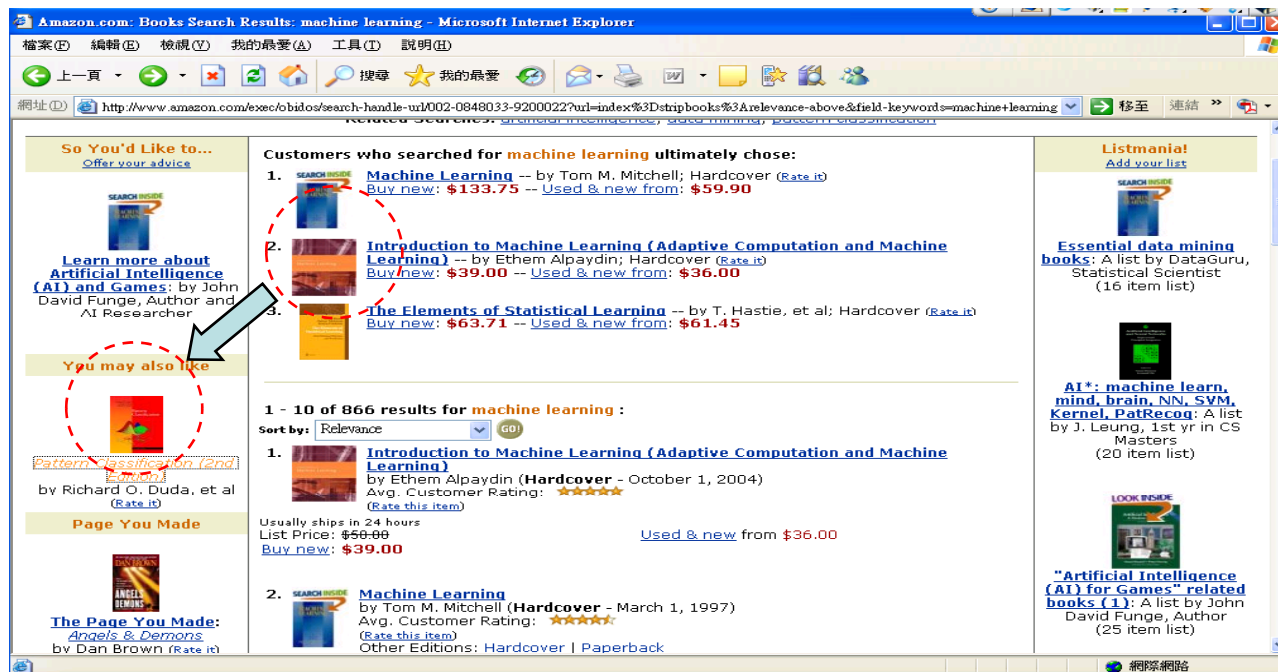| Class Interval (g/gal) | Frequency | Relative Frequency | Density |
|---|---|---|---|
| 1 – < 3 | 12 | 0.193 | 0.0965 |
| 3 – < 5 | 11 | 0.178 | 0.0890 |
| 5 – < 7 | 18 | 0.290 | 0.1450 |
| 7 – < 9 | 9 | 0.146 | 0.0730 |
| 9 – < 11 | 5 | 0.082 | 0.0410 |
| 11 – < 15 | 3 | 0.048 | 0.0120 |
| 15 – < 25 | 4 | 0.063 | 0.0063 |

- Histogram and Frequency table for PM emissions of 62 vehicles driven at high altitude

- Inferential statistics
  - Concerned with the methods that use sample results to help make decisions or predictions about the data (population)
  - Or, the methods that draw conclusions from the data

# Types of Statistics (3/4)

- Example 1
  - A machine makes 1000 steel rods per hour, with a specification of 0.45 $\pm$ 0.02 cm
  - An engineer would like determine the quality/quantity of the production process by randomly draw a sample of rods (say, 50 rods)
  - Given that 92% of the sample meet the specification
    - How likely is the size of difference between the sample proportion and the population proportion?

      Standard derivation (Chapters 2 and 4)
    - How is he confident that the true population proportion will be in 92% $\pm$ *x*%

      Confidence interval (Chapter 5)
    - Can he draw a conclusion that the percentage of good rods is at least 90%

      Hypothesis testing (Chapter 6)
    - ….

# Types of Statistics (4/4)

- Example 2: relationship between two factors/populations



- Association Rule:

  P( buying "Pattern Classification"| buying "Machine Learning" ) = ?

# Popular Software Packages for Statistics

- SPSS
- SAS
- MINITAB
- Microsoft Excel
- …

# Textbook and Reference

- Textbook
  - William C. Navidi, "Statistics for Engineers and Scientists," McGraw-Hill (2 edition, 2007)

- References
  - Prem S. Mann, "Introductory Statistics," Wesley, (6 edition, 2007)
  - D. P. Bertsekas, J. N. Tsitsiklis, "Introduction to Probability," Athena Scientific (2002)

# Topics to be Covered

- Descriptive Statistics (Chapter 1)
- Probability and Common Used Distributions (Chapters 2 & 4, quick review)
- Propagation of Error (Chapter 3)
- Confidence Intervals (Chapter 5)
- Hypothesis Testing (Chapter 6)
- Correlation and Simple Linear Regression (Chapter 7)
- More Topics:
  - Data Analysis and Dimension Reduction
  - Data Cleansing and Presentation
  - Bayesian Decision Theory
  - Parametric Methods - Bias and Variance of the Estimator
  - …

# Grading (Tentatively)

- Midterm and Final: 50%

- Homework: 35%

- Attendance/Other: 15%

- TA: 劉家妏 同學 (碩一)
    - E-mail: acat103@yahoo.com.tw
    - Tel: 29322411ext 208 (資工系208室)