# MIS 985:

# Practical Business Analytics

Course Overview & Introduction

Yihuang K. Kang

# Instructor



- **Yihuang Kang (康藝晃), PhD**
  - Email: ykang@mis.nsysu.edu.tw
  - Office hour: (by appointment)
  - Phone: +886 7-5252000 ext. 4737

- **Shot me an email with "MIS 985 - "**
  - e.g. "MIS 985 - A quick question about Homework #1"
  - Please organize your email before you hit "Send"!

- **Talk is cheap, show me the data**

- **Contact me before "it's too late"**

# Prerequisite

- Basic understanding of relational databases, SQL, data structures, and probability & statistics.

- College-level calculus and matrix operations (Linear Algebra) are required.

- Familiar with at least one high-level programming language. Scientific programming language, such as R, MATLAB, Python, SAS, Julia are preferred.

# Syllabus

- **You will…**
  - exercise logical and computational thinking
  - sharpen your data analytics skills
  - learn how to use R to deal with "big data"
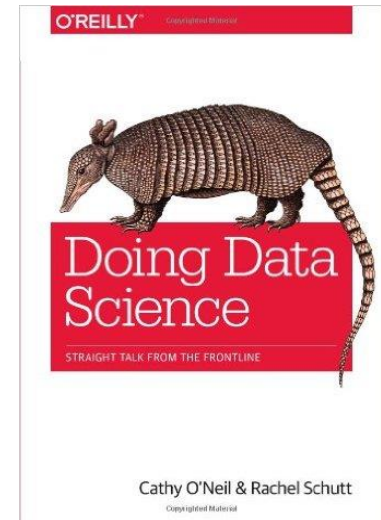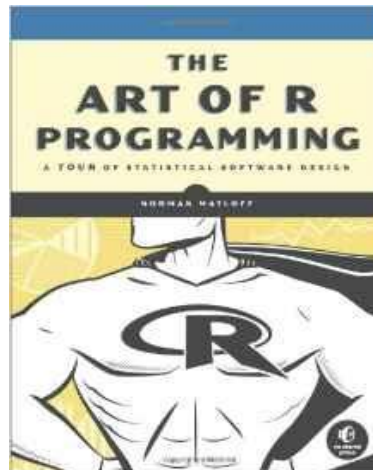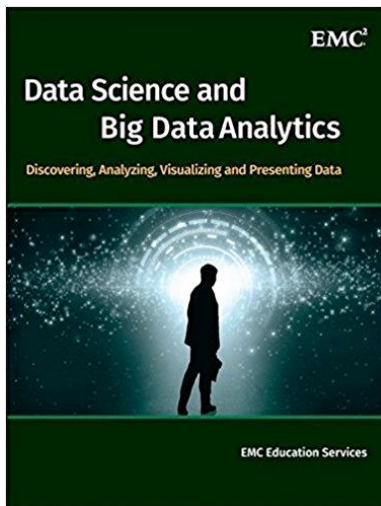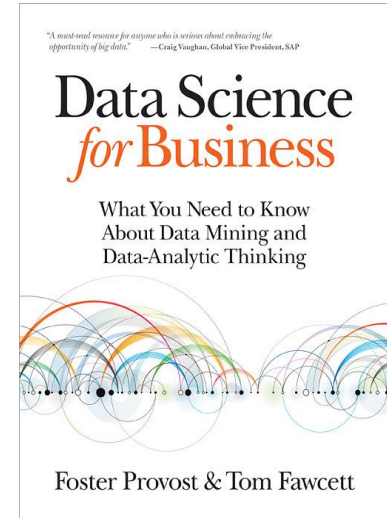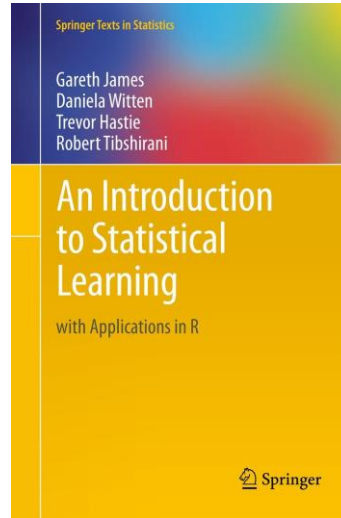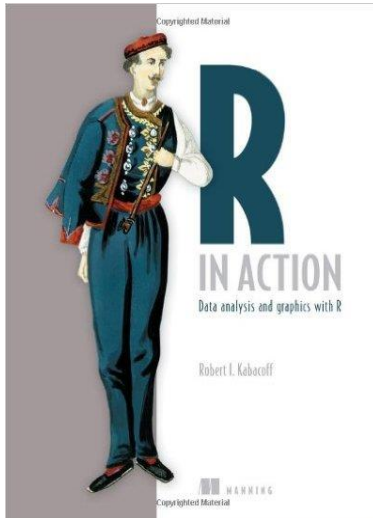
- **Time & Location**
  - Wed 7:10-9pm (online), Sat 1-4pm (monthly on-site, CM 3051),
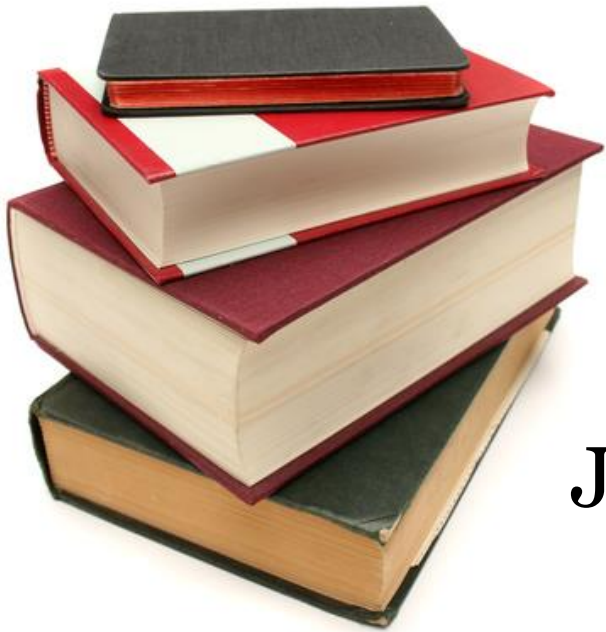
- **Teaching Assistant:**
  - 毛文瑞(rex850327@gmail.com)
  - 黃升泰(rex850327@gmail.com)
  - 周詠捷(johnny83051202@gmail.com)
  - 楊宛蓁(w126251@gmail.com)
  - 郭博文(willy821002@gmail.com)

# Syllabus (cont.)

# Syllabus (cont.)

As for the reading list...

Just read as many as you can!

# Syllabus (cont.)

- **Grading**

  In-class group quiz: 30%          Homework: 30%

  Term project proposal: 20%       Term project defense: 20%

- **Term Project**

  

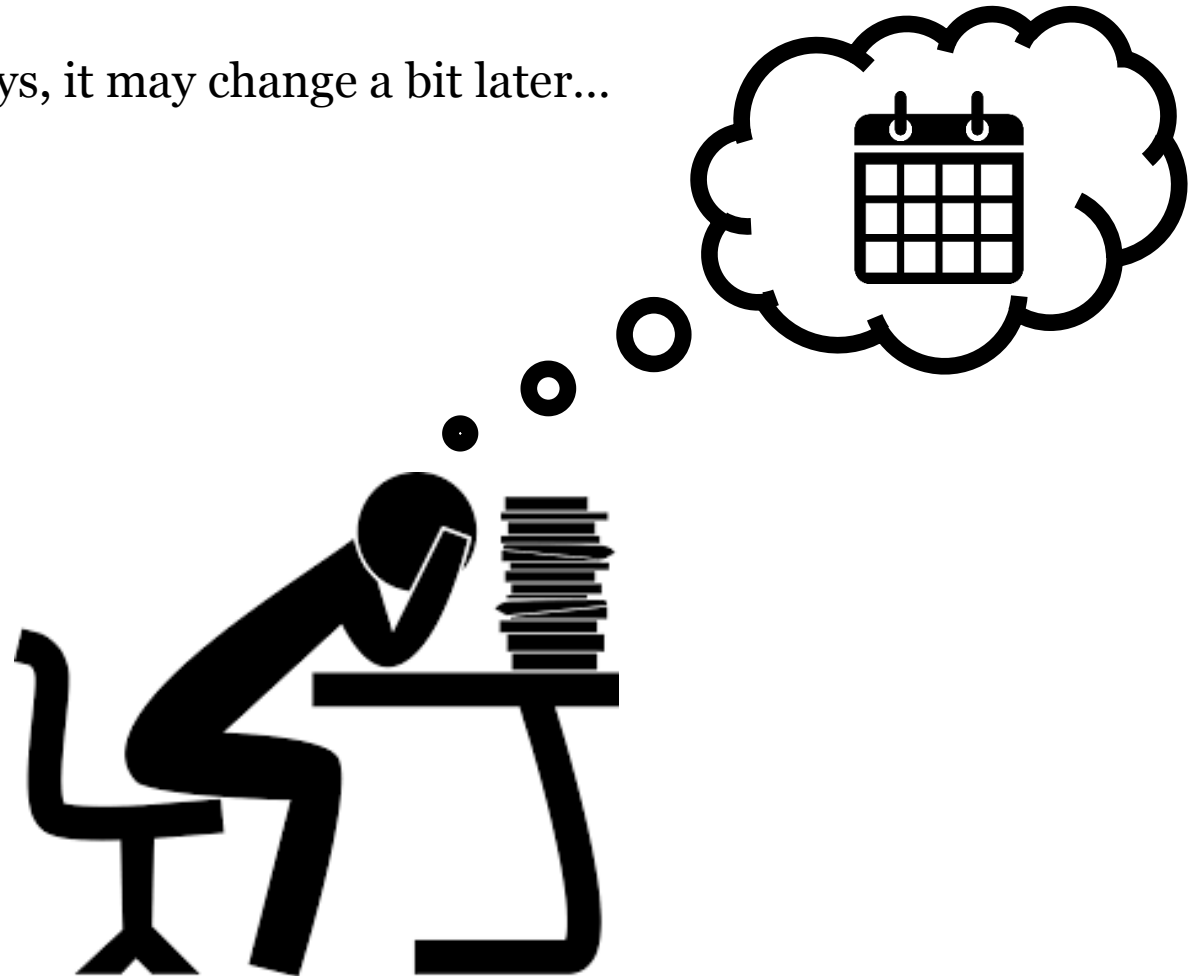  ✓Organize your data science team

  ✓Group of **3-5** people

  ✓Get people with different backgrounds

# Syllabus (cont.)

- **Schedule**

    Well, as always, it may change a bit later...

# How to Survive

✓Practice makes perfect

✓Participate in class discussions

✓Love your data

✓Work ~~hard~~ **smart!**

✓Ask geeks around you

# CM Unified Analytics Platform

- RStudio Server: http://hdp.cm.nsysu.edu.tw:8787/

# CM Unified Analytics Platform(cont.)

# CM Unified Analytics Platform(cont.)

# CM Unified Analytics Platform (cont.)



## II. 系統技術規格：

### A. CPU 平行運算叢集：

- 1 x Master Node: 32 cores, 100G
- 17 x Worker Nodes: 8 cores, 61G
- R 3.4.1, RStudio 1.0.153
- Apache HADOOP 2.7
- Apache Spark 2.1.0
- Apache Zeppelin 0.7.2
    - R 3.4.1
    - Python 2.7.5
    - scala 2.11.8

### B. GPU 深度學習主機：

- 1 x Server: 48 cores, 280G
- 1 x GPU Card: Tesla P100-PCIE-16GB
- R 3.4.1, RStudio 1.0.153
- Python 2.7.5
- MxNet 0.11.0
- Tensorflow-GPU 1.3.0
- Keras 1.2.2 + MxNet 0.11.0
- Keras 2.0.8 + Tensorflow 1.3.0

Please check out http://cm.nsysu.edu.tw/~msrc/wp/ for more information.

# CM Big Data Analytics Software Stack

*"By 2018, the United States will experience a shortage of 190,000 skilled data scientists, and 1.5 million managers and analysts capable of reaping actionable insights from the big data deluge."*

— *McKinsey Report, 2013*

**Big Data, Big Paycheck**

Median salary for analytics professionals and those specifically within data science, by level of experience.

Up to 3 years — Analytics professionals $65,000
Up to 3 years — Data scientists $80,000

4 to 8 years — $85,000
4 to 8 years — $120,000

9+ years — $115,000
9+ years — $150,000

Note: Data do not include managers    Source: Burtch Works    The Wall Street Journal

16

# What is "Data Science"



*Drew Conway's Venn diagram of data science*

# The rise of "Data Scientists"

*"Data Scientist is The Sexiest Job of the 21st Century"*

*—T. Davenport & D.J. Patil, Harvard Business Review*

# The "Dilemma" of the Data Scientists

# Your data careers

# What is "Big Data"

- Many people have defined "Big Data" with 3Vs, 4Vs, 5Vs..., many more Vs!



- My definition is: "*Too much and complicated data to be processed by a single machine with reasonable time or resources*".

# Where does the big data come from?

- **Traditional Data**
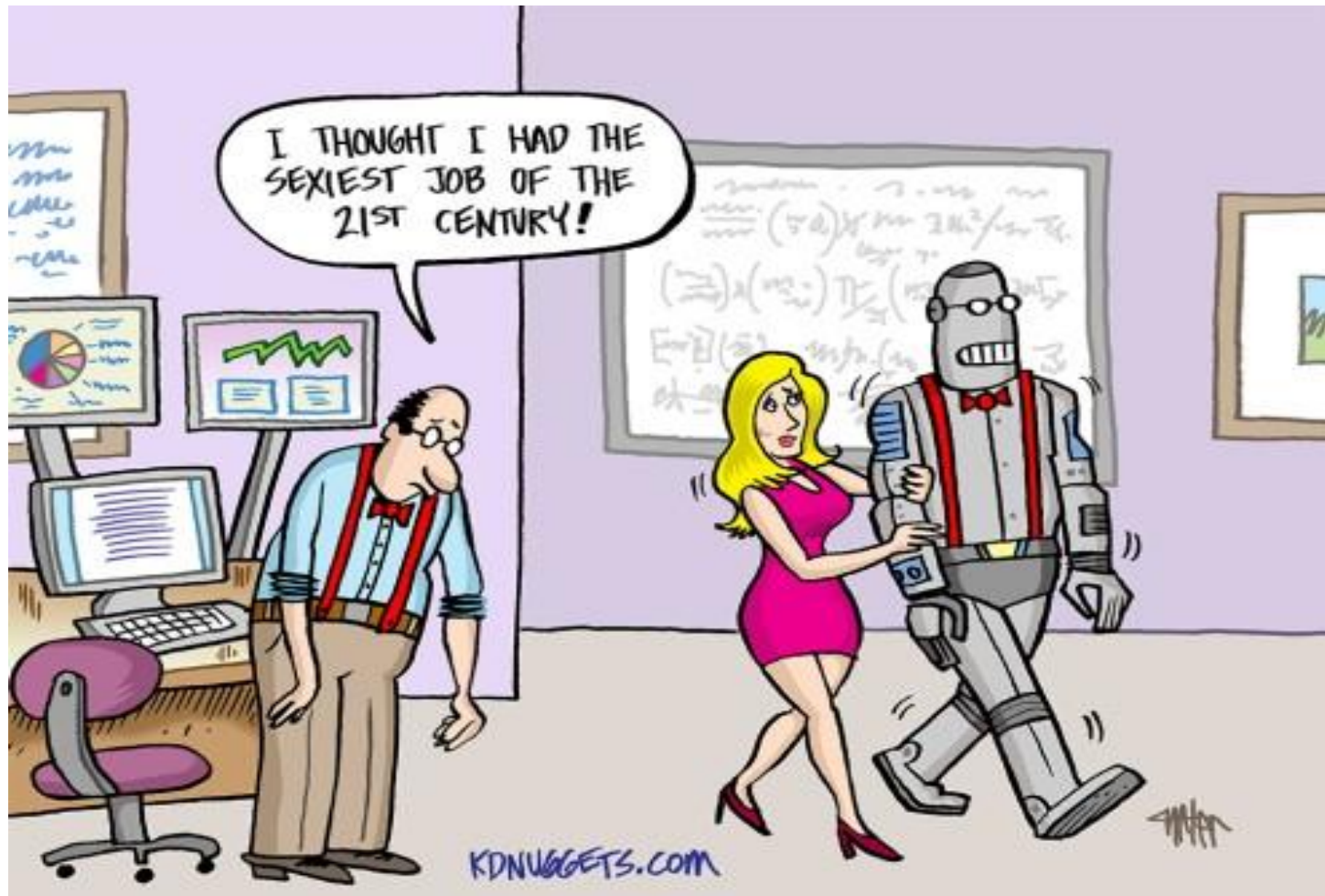  - Any digitized contents and/or archives acquired by traditional ways, e.g. survey data, interview records, and documents.

- **Machine Data**
  - Sensor data, web logs, any log data from monitoring information systems.

- **Network Data**
  - The network of computers (The Internet)
  - The network of people (Social Networks)
  - The network of things (Internet of Things)

# Types of Big Data

- **Structured data**
    - Data with clear schema/metadata/data model that describes & defines how the data elements relate to one another. E.g. relational databases, data cubes/warehouses.
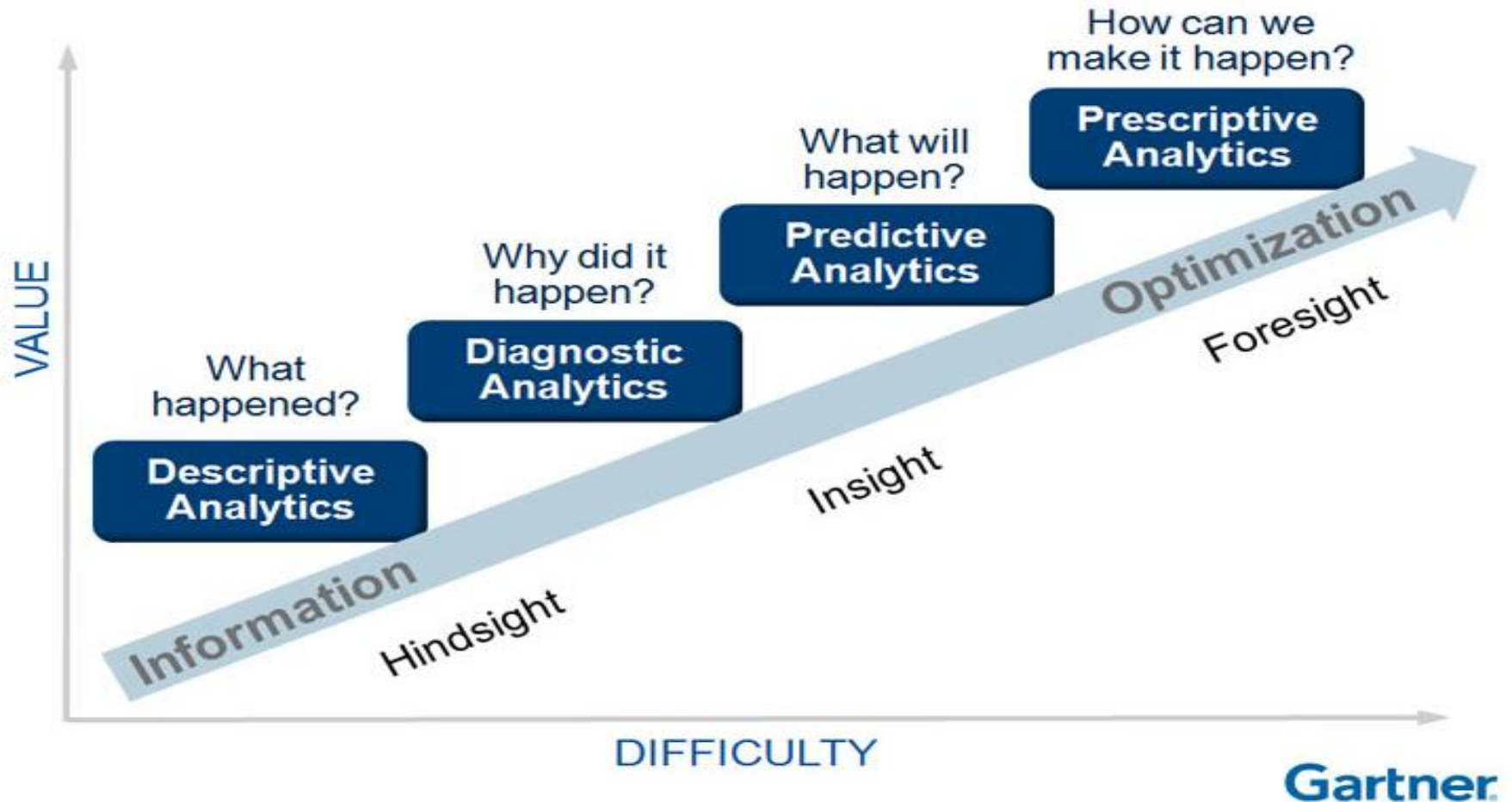
- **Semi-Structured data**
    - Data with only tag/field definitions but without formal structures of data models to define relations. E.g. data used in information exchanges, such as XML & JSON. Emails/pictures/other files with tags/field definitions.

- **Unstructured data**
    - Unorganized data without any pre-defined schema. E.g. body of an e-mail message, pictures, audio, and video.

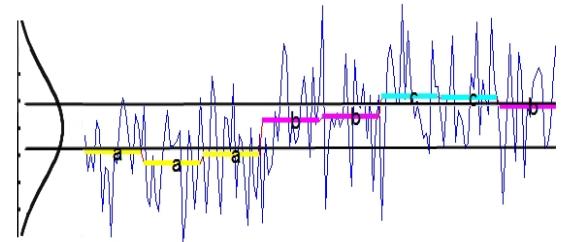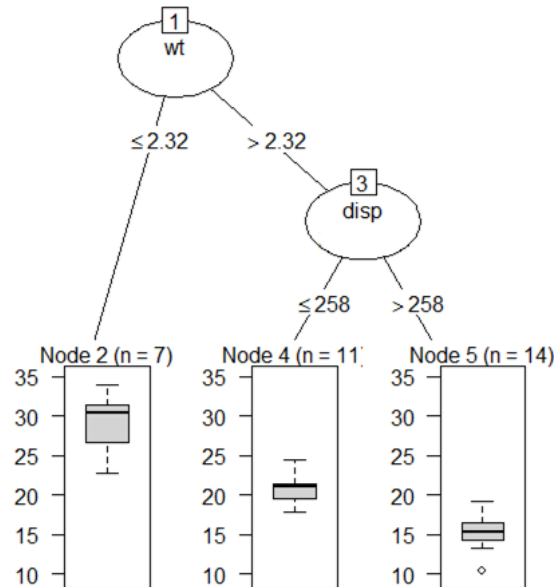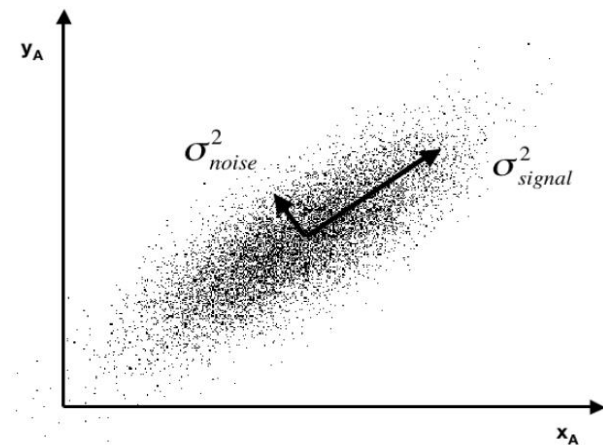# Wait.. we're talking about "Big Data Analytics"?

# Data is NOT always the cure!

- The "Big Data" does eliminate intuition. However, our interpretations of it have great impact on the results. Let's check out [this article in New York Time](). It says "*Let's put everything in and let the data speak for itself.*" This is a bit horrible quote and don't let it mislead you.

"*...Data is just a quantitative, pale echo of the events of our society...*".

–*O'Neil, "On Being a Data Skeptic"*

# Statistical Machine Learning

# MapReduce Design Patterns

- We will surely do more than just the "word count"!

**Filtering**



**Aggregation (crosstab)**



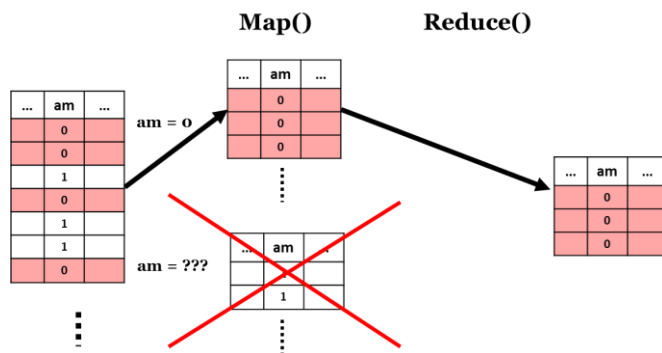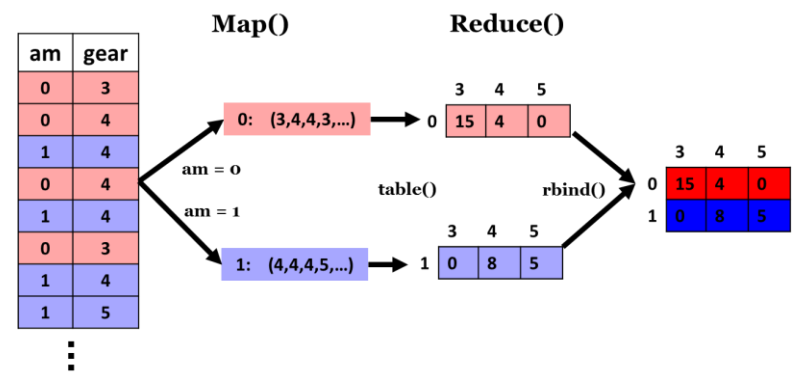**Split-Apply-Combine (e.g. model fitting)**



**Sorting patterns,  Join patterns, ….**

**And many more!**

# High-performance R Programming

- We still don't know much about R's own limitations and capabilities when coping with Big Data. Why my R code is so slow? How to evaluate my R code?

- We will be discussing vectorized and functional programming, and why they matter in the age of Big Data.

- We will also be discussing how to tweak programs by writing more functional, primitive, and parallel R code, as well as how to use more CPU cores on both a single and a cluster of machines!

# Two Ecosystems of Big Data Analytics

**In-database Analytics**



**Table Partitions**

**Massive Parallel Processing Databases**
**(e.g. Teradata, PostgreSQL, Greenplum,... etc.)**

**VS.**

**In-memory Analytics**

**Distributed DataFrames/Datasets**

*"AI is the New Electricity."*

— *Andrew Ng (吳恩達)*

# What would make an AI company?

- Strategic Data Acquisition

- Centralized Data Warehouse (Unified Data Analytics Platform)

- Pervasive Automation

# Trends in Big Data Analytics

- ✓ The flood of "data lake"

- ✓ The rise of out-of-memory analytics/algorithms

- ✓ The dawn of *fast scalable data applications*

- ✓ The use of *in-memory* & *in-database* computations

- ✓ The pursuit of accountable and transparent AI

- ✓ The fall of business without AI strategies

# Your homework this week

❑ Make a choice. Quit or stay. We will be getting an account for the access to CM Big Data Analytics Platform next week.

❑ Review R programming. Especially those of you who are not familiar with any scientific computing languages.

❑ Get the textbooks & papers and start reading!

# See you next week!