

Lecture 1 : Introduction

§1 Data science

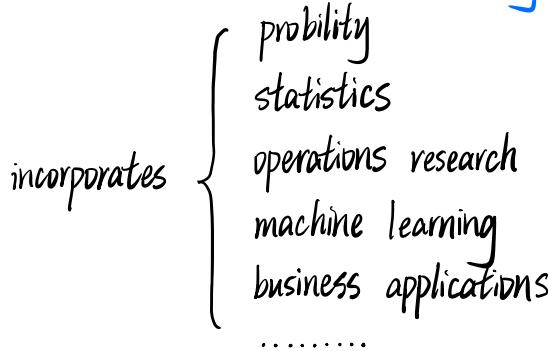
1. What's data science

1^o deal with vast volumes of data

2^o using modern tools and techniques

3^o to find unseen patterns / derive meaningful information / make business decisions.

2. Data science is not one single subject!



3. What do you need

1^o Analytical skills: solid mathematical / technical foundation

2^o Business understanding: industrial expertise

4. Some examples



Blind Box Example

- Ten of you went to a blind box seller for the following blind box. You were told the chance to get a Harry Potter versus a Lord Voldemort is fifty-fifty.
- You all wish to get a Harry Potter.



Blind Box Example

- If the seller told the truth, the probability that all ten of you got a Voldemort is < 0.001 . \rightarrow Probability
- Given the above observation, you are pretty confident that the seller was lying. What's the true chance of getting a Harry Potter versus Voldemort? \rightarrow Statistics
- Suppose you are a honest seller. How should you design the blind box and its price to maximize your own profit? \rightarrow Optimization
- Design a Cartoon for Harry Potter. Which one is more likely to be popular?

• Jack opened a shoe store and he didn't want to miss this opportunity.

• He plans to do some publicity, in order to increase the exposure and attract more customers to buy his products.

• There are many advertising channels.



- How much shall Jack invest in each channel?
- Do you have any advices?



- Collect other people's revenue performance with different investments.
- Let investment be X , and the revenue be R .
- Data: $(X_i, R_i), \dots$
- From data, find the relation between X and R .
- Given the pattern as $R(X)$, Jack's decision shall be

$$\operatorname{argmax}_X R(X)$$

The value of X such that $R(X)$ is maximized ³²

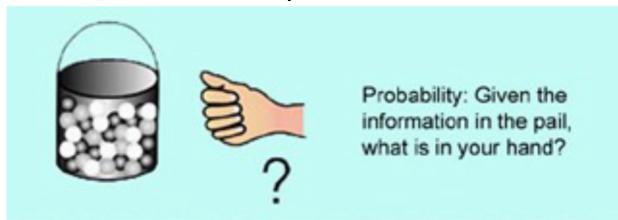
§2 Key steps in data science

1. Programming (编程) (How to collect the past data?)

1^o creating instructions that tell a computer how to perform a task

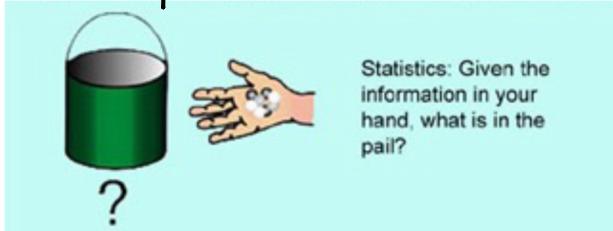
2^o fundamental part

2. Probability (概率) theory (How to describe the data pattern?)
- 1º A formality to make sense of the world in terms of uncertainty
 - 2º Numerical descriptions of how likely an event is to occur

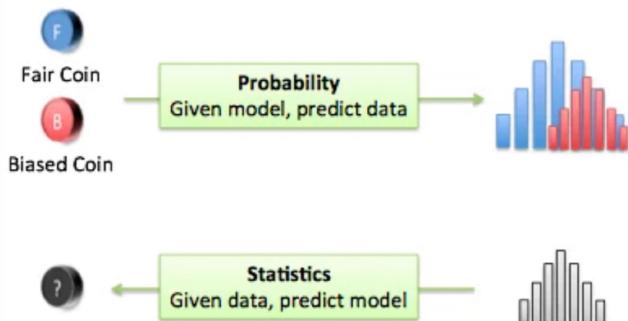


Statistics (统计学) (Is the data pattern we found reliable?)

- 1º The discipline that extracts correct information from data.



* Probability versus statistics



3. Simulation (模拟) (Given pattern, can we calculate objective for different strategy?)

1º a computerized mathematical technique

2º generate random samples

3º aid quantitative analysis and decision making

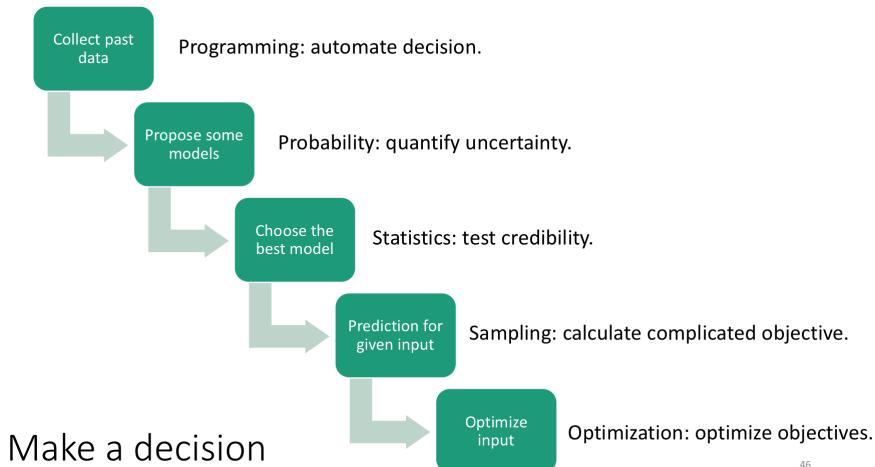
4. Optimization (最优化) (What decision should be made to maximize objective?)

1º The selection of a best element, with regard to some criterion, from some set of available alternatives

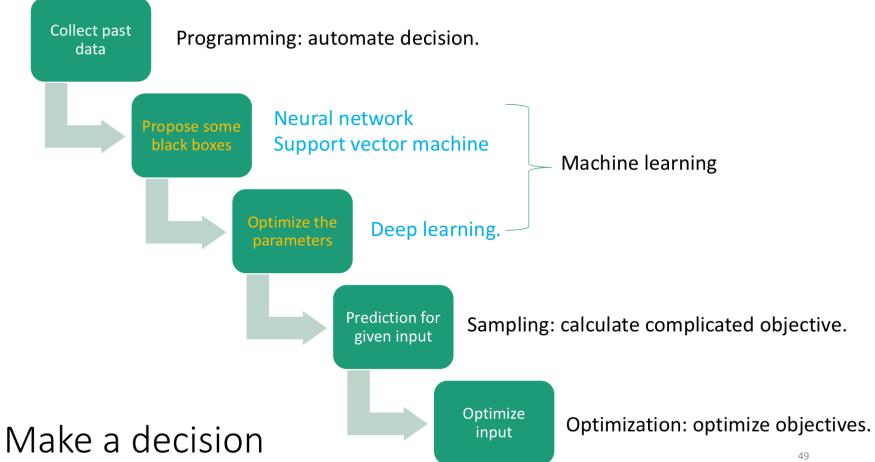
$$\min_{x \in \mathbb{R}} (x^2 + 1)$$

- 5* Machine learning (机器学习) (If step 2 is too complicated)

Allow software applications to become more accurate at predicting outcomes without being explicitly programmed to do so.



46



49

Questions

- Given the data, you want to predict how the data is generated. Which theory you may use?
 - Probability
 - Statistics
- Given the profit as a function of strategies, you want to maximize the profit. Which theory you may use?
 - Simulation
 - Optimization

§3 Why study data science

1. Data science can be used in...

- 1º retail (零售)
- 2º healthcare
- 3º finance
- 4º education
- 5º transportation
- 6º agriculture
- 7º startups (创业)
- 8º pharmaceutical industry (医药产业)
- 9º art industry