**DDA2001: Introduction to Data Science**

# Introduction

**2023/01/08**

**Zicheng Wang**

# An appetizer of data science application

The Math Behind Basketball's Wildest Moves | Rajiv Maheswaran | TED Talks

Data science aims to understand the world by collecting and analyzing data, and makes better decisions for the future.

# Outline

Course organization

Course introduction

# Course organization

# Course staff

- Instructor:

  - *Zicheng Wang*

    - Office:          DY 322B, TXC711
    - Office Hours:    Monday 9am – 10am
    - Email:          wangzicheng@cuhk.edu.cn

- Lectures:

    - Room:           Teaching Complex C202

    - Times:          Mon Wed 10:30 – 11:50am

# Course staff (TAs)

| | |
|---|---|
| •Chenguang Wang | 王晨光 |
| Email: | 222043009@link.cuhk.edu.cn |
| •Chao Yang | 杨超 |
| Email: | 222043011@link.cuhk.edu.cn |
| •Weihuang Wen | 温伟煌 |
| Email: | 223040256@link.cuhk.edu.cn |
| •Sheng Xu | 徐圣 |
| Email: | 223040246@link.cuhk.edu.cn |
| •Wenjun Zheng | 郑文军 |
| Email: | 118010446@link.cuhk.edu.cn |

# Schedule

| | Monday | Tuesday | Wednesday | Thursday | Friday |
|---|---|---|---|---|---|
| 8:30--9:00 | | | | | |
| 9:00--9:30 | | | | | |
| 9:30--10:00 | | | | | |
| 10:00--10:30 | | | | | |
| 10:30--11:00 | Lecture | | Lecture | | |
| 11:00--11:30 | Zicheng Wang | | Zicheng Wang | | |
| 11:30--12:00 | Teaching Complex C202 | | Teaching Complex C202 | | |
| 12:00--12:30 | | | | | |
| 12:30--13:00 | | | | | |
| 13:00--13:30 | | | | | |
| 13:30--14:00 | Lecture | Lecture | Lecture | Lecture | |
| 14:00--14:30 | Shuang Li | Pin Gao | Shuang Li | Pin Gao | |
| 14:30--15:00 | ZhiXin 111 | Teaching Complex C303 | ZhiXin 111 | Teaching Complex C303 | |
| 15:00--15:30 | | | | | |
| 15:30--16:00 | Lecture | | Lecture | | |
| 16:00--16:30 | Jiaqi Lu | | Jiaqi Lu | | |
| 16:30--17:00 | Teaching Complex B206 | | Teaching Complex B206 | | |
| 17:00--17:30 | | | | | |
| 17:30--18:00 | | | | | |
| 18:00--18:30 | | Tutorial 03 | Tutorial 05 | | |
| 18:30--19:00 | | TA 302 | TA 310 | | |
| 19:00--19:30 | Tutorial 01 | Tutorial 04 | Tutorial 06 | | |
| 19:30--20:00 | Teaching Complex B205 | TA 302 | TA 310 | | |
| 20:00--20:30 | Tutorial 02 | | Tutorial 07 | | |
| 20:30--21:00 | Teaching Complex B205 | | TA 310 | | |

# Tutorials

- **Tutorials start on the <u>second</u> week**

| | Tutorial Session | TA | Location | Time | |
|---|---|---|---|---|---|
| **Probability** | T01 | Chenguang Wang | TCB 205 | Mon:19:00-19:50 | |
| | T02 | Chenguang Wang | TCB 205 | Mon:20:00-20:50 | |
| | T03 | Chenguang Wang | TA302 | Tues: 18:00-18.50 | |
| | T04 | Chenguang Wang | TA302 | Tues: 19:00-19.50 | |
| | T05 | Chenguang Wang | TA310 | Wed: 18:00-18:50 | |
| | T06 | Chenguang Wang | TA310 | Wed: 19:00-19:50 | |
| | T07 | Chenguang Wang | TA310 | Wed: 20:00-20:50 | |
| **Statistics** | Tutorial Session | TA | Location | Time | |
| | T01 | Sheng Xu | TCB 205 | Mon:19:00-19:50 | |
| | T02 | Sheng Xu | TCB 205 | Mon:20:00-20:50 | |
| | T03 | Sheng Xu | TA302 | Tues: 18:00-18.50 | |
| | T04 | Sheng Xu | TA302 | Tues: 19:00-19.50 | |
| | T05 | Sheng Xu | TA310 | Wed: 18:00-18:50 | |
| | T06 | Sheng Xu | TA310 | Wed: 19:00-19:50 | |
| | T07 | Sheng Xu | TA310 | Wed: 20:00-20:50 | |
| **Optimization** | Tutorial Session | TA | Location | Time | |
| | T01 | Wenjun Zheng | TCB 205 | Mon:19:00-19:50 | |
| | T02 | Wenjun Zheng | TCB 205 | Mon:20:00-20:50 | |
| | T03 | Wenjun Zheng | TA302 | Tues: 18:00-18.50 | |
| | T04 | Wenjun Zheng | TA302 | Tues: 19:00-19.50 | |
| | T05 | Wenjun Zheng | TA310 | Wed: 18:00-18:50 | |
| | T06 | Wenjun Zheng | TA310 | Wed: 19:00-19:50 | |
| | T07 | Wenjun Zheng | TA310 | Wed: 20:00-20:50 | |
| **ML** | Tutorial Session | TA | Location | Time | |
| | T01 | Weihuang Wen | TCB 205 | Mon:19:00-19:50 | |
| | T02 | Weihuang Wen | TCB 205 | Mon:20:00-20:50 | |
| | T03 | Weihuang Wen | TA302 | Tues: 18:00-18.50 | |
| | T04 | Weihuang Wen | TA302 | Tues: 19:00-19.50 | |
| | T05 | Weihuang Wen | TA310 | Wed: 18:00-18:50 | |
| | T06 | Weihuang Wen | TA310 | Wed: 19:00-19:50 | |
| | T07 | Weihuang Wen | TA310 | Wed: 20:00-20:50 | |

# Tentative Calendar

| Weeks | Content/ topic/ activity |
|---|---|
| | **Course Introduction** |
| 1 | Course plan, grading, overview, examples, and applications, etc. |
| | **Probability** |
| 3 | 1 Terminologies: sample spaces, probability function, event, etc.. <br> 2. Main Properties: mean, variance, CDF. <br> 3. Common Discrete RV: Bernoulli, Binomial, Poisson, Geometric. <br> 4. Common Continuous RV: uniform, normal, exponential. <br> 5*. Advanced Concepts: correlation and conditional probability (Bayes rule). <br> 6#. Additional examples and exercises. |
| | **Statistics** |
| 3 | 1. Sampling: Monte Carlo (integration, MDP, etc.). <br> 2. Sampling: generating uncommon RV. <br> 3. Point estimator: MLE. <br> 4. MLE Application: Regression Analysis. <br> 5*. Confidence interval <br> 6#. Additional examples and exercises. |
| 0.5 | **Midterm Review** |

# Tentative Calendar

| | |
|---|---|
| | **Optimization** |
| 3.5 | 1. Problem formulation (objective function, constraints, decision variables, etc.) Application examples in revenue management, online advertising, etc.<br>2. Optimization: convex set.<br>3. Convex Optimization: convex function.<br>4. Convex Optimization: convex function Property.<br>5. Convex Optimization: quasi-convex function.<br>6*. Gradient Descent and Stochastic Gradient Descent<br>7#. Additional examples and exercises. |
| | **Machine Learning** |
| 2.5 | 1*. Supervised learning:  K-Nearest Neighbors (K-NN).<br>2*. Supervised learning: Logistic regression, Neural Nets.<br>3*. Unsupervised learning: Clustering (K means).<br>4**. Cross Validation.<br>5**: Deep Learning and Reinforcement Learning. |
| 0.5 | **Final Review** |

\# Problem exercise.
\* Only concepts will be examined.
\*\* Contents that will be covered if time permits. If covered, only concepts will be examined.

# Learning Assessment Strategy

| Component/ method | % weight |
|---|---|
| Assignments | 15% |
| Quizzes | 15% |
| Mid-Term Exam | 30% |
| Final Exam | 40% |

- **3 to 5 Guest Lectures on Fridays (Dates TBD). Quizzes on Guest Lectures.**

- **No late homework will be accepted.**

- **Midterm (tentative): TBD**

- **There will be no make-up midterm exams**

- **Zero tolerance policy when it comes to cheating so, DON'T CHEAT**

- **Missing the midterm or final exam without prior notification to and approval of the instructors will automatically result in the "0" grade for the exam.**

# About me

- **My academic pathway**
  - Undergraduate: Mathematics, University of Minnesota
  - Graduate: PhD in Industrial and Systems Engineering, University of Minnesota (advisor: Kevin Leder)
  - Join SDS of CUHK(SZ) since July 2023, and currently an assistant professor
- **My Research Interest related to "Data Science": Operations Research**
  - Applied probability and its application in mathematical modeling of cancer and operations management
  - Delineate tumor dynamics
  - Analyze complex spatial systems

# Course introduction

## C. F. Jeff Wu

Carver Professorship at the University of Michigan. ... He popularized the term "data science" and advocated that statistics be renamed data science and ...

## Statistics Data Science

Statistics Data Science ? = C. F. Jeff Wu. University of Michigan, Ann Arbor. What is "Statistics"? •. A Statistical Trilogy. Frontier and Beyond. • A Bold ...

13 页

---

## 吴建福 (C. F. Jeff Wu)

统计师

吴建福，美國台裔統計學家，喬治亞理工學院工業及系統工程系可口可樂講座教授。主要從事工業統計與質量科學的研究及應用。研究成果包括EM算法收斂性的討論。 维基百科

**出生地：** 台灣新竹

**教育背景：** 加利福尼亞大學， 國立臺灣大學

**知名校友：** 蒂爾桑卡·達斯古普塔， Jiahua Chen， Jun Shao， 更多

**学术顾问：** 彼得·畢克爾

**奖项：** 費希爾獎， Shewhart Medal

# What is data science

- Data science is the domain of study that deals with vast volumes of data using modern tools and techniques to find unseen patterns, derive meaningful information, and make business decisions.

- Data science is not one single subject! It incorporates probability, statistics, operations research, machine learning, and business applications…

- You need a solid mathematical/technical foundation, together with some industrial expertise in order to be a data scientist.

- Analytical skills + business understanding.

# How are they related?

Probability

Statistics

Optimization

Machine learning

GLOBAL SHOPPING FESTIVAL 2021

The festival has just passed.
Did you buy a lot of things?

剁手人、吃土人、尾款人、晚八人

**1.剁手人 (Duò shǒu rén)**: Literally translates to "hand-chopping people". It's a humorous way to describe someone who can't resist buying things online, implying they might as well chop off their hands to stop spending.

**2.吃土人 (Chī tǔ rén)**: Literally means "eating dirt people". This phrase is used to describe people who have spent so much money, particularly on shopping, that they are left with no money for essentials, metaphorically left to 'eat dirt'.

**3.尾款人 (Wěi kuǎn rén)**: This translates to "tail payment people". It's used to describe people who make a down payment for a product (often during a sale like Singles' Day in China) and are waiting to pay the balance amount.

**4.晚八人 (Wǎn bā rén)**: This phrase is a bit more context-dependent. Literally, it means "late-eight people". This could refer to people who work late (until 8 PM), or it could be a specific cultural reference that needs more context to accurately translate.

- To stand out from the event, various promotional rules were designed



**PRICING STRATEGIES IN MARKETING**

Bundle pricing strategy

Discount pricing strategy

Economy pricing strategy

Freemium pricing strategy

Geographic pricing strategy

Premium pricing strategy

Psychological pricing strategy

Penetration pricing strategy

www.financialfalconet.com

# Things to consider when choosing a pricing strategy

**Cost of Products**

**Customers**

**Market Positioning**

**Competitors**

**Profit Margins**

**Geography**

**Product Offerings**

E-commerce merchant

Forced to become mathematicians

https://shanshu.ai/news/detail?id=243

Do I have to set the prices for all these?



这些都要我来定价格？
听我说，谢谢你

- Even if you are already good at these pricing techniques, it is still very difficult to implement them.

# How complicated are the price rules?

1, 2, 3, click the link, one yuan for three bottles of craft beer

Has the price not changed?

Sorry!

# How complicated are the price rules?

Sorry, setting this price was too troublesome. Full refunds for all purchases and free giveaways! Let's try another round!

对不起，这个价格设置太麻烦了 已购买全部退款免费送！再抢一轮

哈？还是原价，又翻车了？

对不起，我是新人主播

What? Still the original price? Crashed again?

Sorry, I am a new broadcaster.

https://shanshu.ai/news/detail?id=243

# How complicated are the price rules?

There may be hundreds of SKUs (stock-keeping units) in one live broadcast

For e-commerce platforms, what they are facing can be more than millions of SKUs

https://shanshu.ai/news/detail?id=243

Say NO to manual pricing!

How to design prices more
**scientifically and efficiently**?

1. Let a machine think like you

2. Let the machine automatically design prices

# Data science

- Let a machine think like you
  - Extract knowledge from **data** sets

- Let the machine automatically design prices
  - Apply the knowledge and actionable insights from data to solve problems in a wide range of application domains

# To be a good data scientist,

- Programming languages + analytical skills + business understanding.

- It incorporates probability, statistics, operations research, machine learning, and business applications…

# A concrete example

How to price a flight ticket?

# 1. Data collection? **Programming**.

- Obtain historical data, obtain current competitor prices, etc.

1. Data collection? **Programming**.

- Obtain historical data, obtain current competitor prices, etc.

- Programming is the process of creating a set of instructions that tell a computer how to perform a task

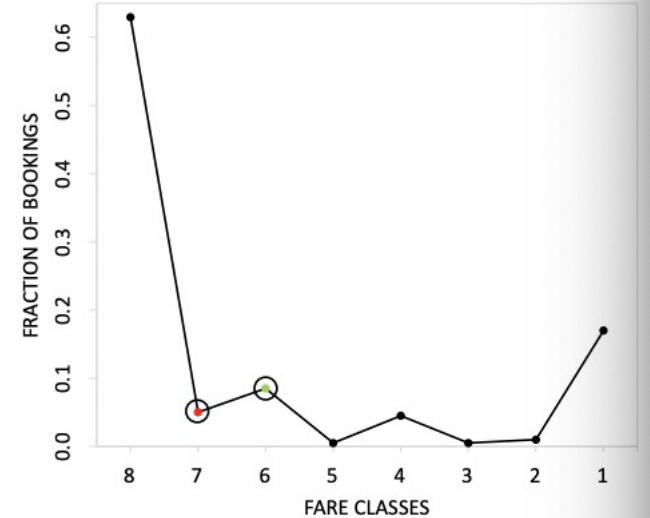- Learning to program is to free your hands

- Programming is a fundamental part of data science.

- Programming Languages for Data Science: python, R, C, C++, …

- In this course,
  - We will give some coding demos
  - There will be no coding questions in the homework and the exams.

1. Data collection? **Programming**.
2. How to describe the data pattern? Probability theory.

**You observe:**

The product with the lowest price always sells significantly higher than other products.

How to describe the data generation process?

Daniel Little McFadden (born July 29, 1937) is an American econometrician who shared the 2000 Nobel Memorial Prize in Economic Sciences with James Heckman. McFadden's share of the prize was "for his development of theory and methods for analyzing discrete choice".
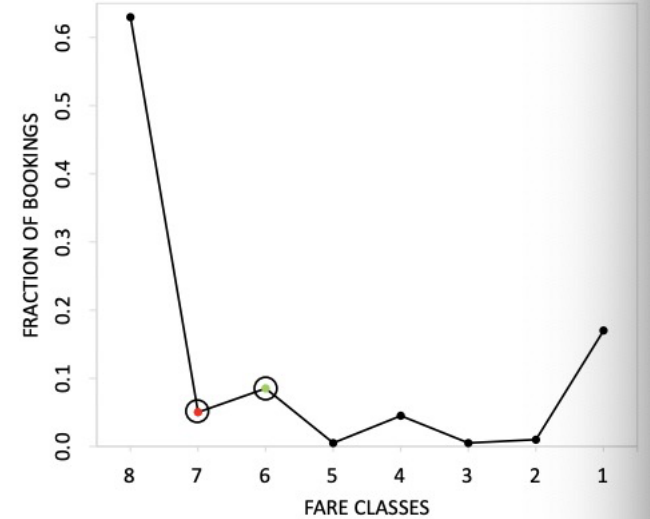
Daniel McFadden - Wikipedia

1. Data collection? **Programming**.
2. How to describe the data pattern? Probability theory.

Consider a toy model:

With a set of products S, assume that the purchase probability of product i

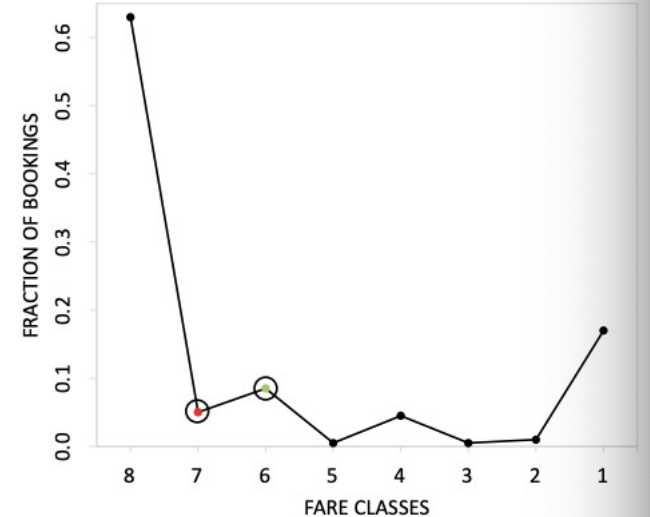$$P_i(S) = \frac{a_i + x_i v_i}{1 + \sum_{j \in S}(a_j + x_j v_j)}$$

- $a_i$ and $v_i$ are known functions of prices
- $x_i \in \{0,1\}$ and $x_i=1$ if and only if this product has the lowest price



FRACTION OF BOOKINGS

FARE CLASSES

1. Data collection? **Programming**.
2. How to describe the data pattern? Probability theory.
3. How to choose a reliable data pattern?  Statistics.

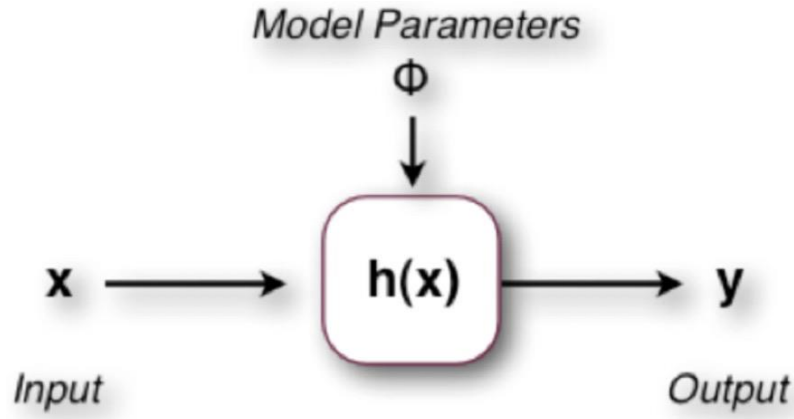$$P_i(S) = \frac{a_i + x_i v_i}{1 + \sum_{j \in S}(a_j + x_j v_j)}$$

- How to determine the values of $a_i$ and $v_i$?
- Statistics can be used to test the quality of the model.

1. Data collection? **Programming**.
2. How to describe the data pattern? Probability theory.
3. How to choose a reliable data pattern?  Statistics.
4. How to learn a complicated function? Machine Learning.

$$P_i(S) = \frac{a_i + x_i v_i}{1 + \sum_{j \in S}(a_j + x_j v_j)}$$

- $a_i$ and $v_i$ may depend on product attributes and consumer profiles.

1. Data collection? **Programming**.
2. How to describe the data pattern? Probability theory.
3. How to choose a reliable data pattern?  Statistics.
4. How to learn a complicated function? Machine Learning.
5. What decision should be made to maximize objective? Optimization

$$\max_r \sum_{j \in S} r_i P_i(S|r)$$

- Choose a price vector to maximize the above objective.

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

1. Data collection? **Programming**.
2. How to describe the data pattern? Probability theory.
3. How to choose a reliable data pattern?  Statistics.
4. How to learn a complicated function? Machine Learning.
5. What decision should be made to maximize objective? Optimization
6. Calculate a complicated objective? Simulations

$$\sum_{j \in S} r_i P_i(S|r)$$

A computerized mathematical technique that allows people to generate random samples and aid quantitative analysis and decision making.

Collect past data

Programming: automate decision.

Propose some models

Probability: quantify uncertainty.

Choose the best model

Statistics: test credibility.

Prediction for given input

Sampling: calculate complicated objective.

Optimize input

Optimization: optimize objectives.

Make a decision

Collect past data

Programming: automate decision.

No models to propose?

Propose some models

Probability: quantify uncertainty.

Choose the best model

Statistics: test credibility.

Prediction for given input

Sampling: calculate complicated objective.

Optimize input

Optimization: optimize objectives.

Make a decision

47

Collect past data

Programming: automate decision.

Propose some black boxes

Neural network
Support vector machine

Optimize the parameters

Deep learning.

Prediction for given input

Sampling: calculate complicated objective.

Optimize input

Optimization: optimize objectives.

Make a decision

Collect past data

Programming: automate decision.

Propose some black boxes

Neural network
Support vector machine

Machine learning

Optimize the parameters

Deep learning.

Prediction for given input

Sampling: calculate complicated objective.

Optimize input

Optimization: optimize objectives.

Make a decision

# Distinction with other courses

- This is an interdisciplinary course and will cover broadly
- More focus on the big picture and the connection between different concepts
- Probability: STA2001/STA2003
- Statistics: STA2002/STA2004
- Optimization: MAT3007
- Machine Learning: DDA3020