### **Data Science & AI for Economists**

Lecture 0: Introduction

Zhaopeng Qu Business School,Nanjing University August 27 2025

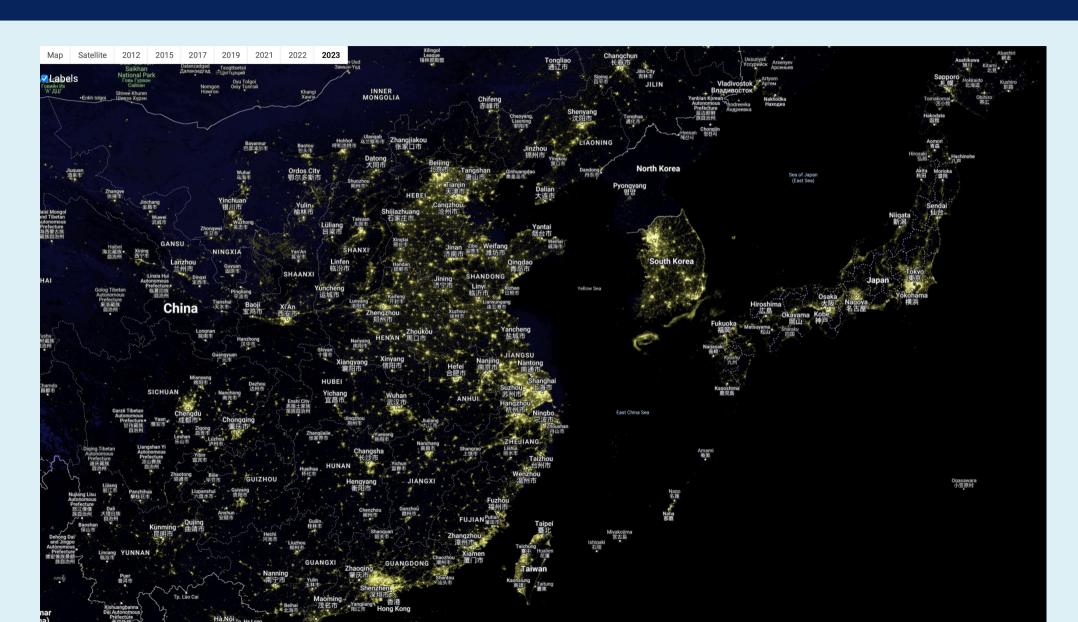


### Introduction to Data Science and AI

### Case #1: Alternative economic indicators

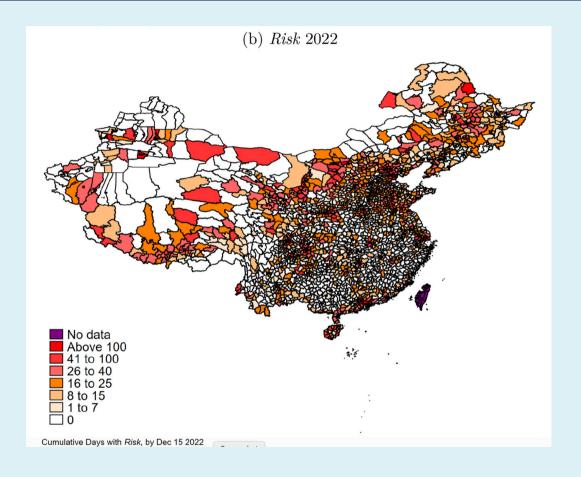
- Question: How do we measure a country's economic performance without published or convincing statistic, like North Korea?
- Answer: Use alternative indicators, such as satellite images of nighttime lights, to estimate economic activity.

### Case #1: Alternative economic indicators



## Case #2: Measuring the intensity of a policy

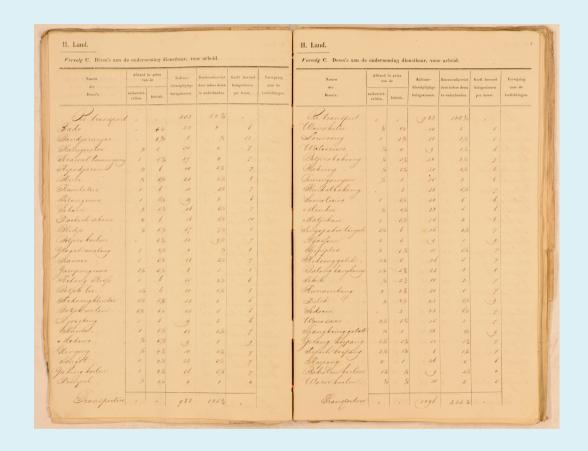
- How to measure the intensity of a policy, such as the COVID-19 lockdown policy in China?
- As is well known, there is huge difference for the intensity of the COVID-19 lockdown policy in China across regions and time.



Gong et al(2024)

## Cases #3: Collecting information from documents

- Information (data) recorded for a long time in terms of documents, papers, books and other forms.
- Traditional way to collect data from old documents is by-hands.



Dell(2020)

## Cases #4: Big data and its applications: Job Postings

- Job postings are a rich source of information about the labor market.
- The observations are huge over **100 millions** job postings from 2013 to 2023.
- The data is not structured, and the format is not consistent, with a lot of missing values.
- More important, the file size is huge, over 300GB.

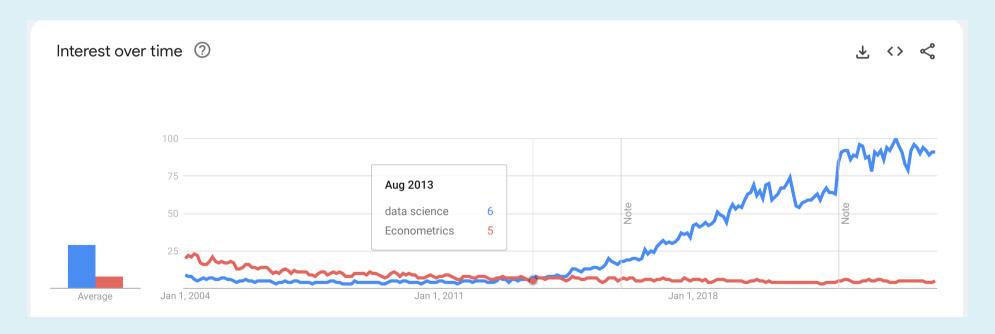


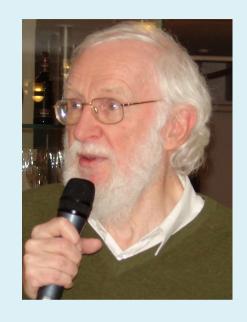
Zhu(2023)

### New Questions need new tools and skills

- New data sources and neew data types are emerging, which make social science research more challenging and exciting.
- It requires new tools and skills to obtain, process, analyze and visualize various data.
- It is the job of **Data Science**.

- Data Science is an **interdisciplinary** academic field that compromises mainly concepts from statistics, computer science, and information science, etc.
- Although scientist have always relied on data to test their theories and make predictions, the term "Data Science" has only recently become popular during the last decade.





Peter Naur(1928-2016)

- Danish computer scientist, Turing Award winner (2005)
- Firstly coined the term "Data Science" in 1960s.



William S. Cleveland(1943-)

- American computer scientist and statistician,
- Formally defining and naming the field of Data Science in 2001.



Joshua Angrist(1960-)

- American Economist and Professor at MIT
- 2021 Nobel Prize co-winner

"One thing I always say is(that) econometrics is the original data science. Before there was data science, there was econometrics." in an public interview at 2021 on Youtube.

• And his econometrics course in MIT now are renamed as **Econometrics Data Science**.



Hadley Wickham(Posit)

- Committee of Presidents of Statistical Societies(COPSS) Presidents' Award winner(2019).
- Chief Scientist at Posit(Former named as RStudio).
- One of Leading figure in the field of Data Science with some most popular packages in R.
  - ggplot2 and tidyverse.

"Data science is an exciting discipline that allows you to transform raw data into understanding, insight, and knowledge." in *R for Data Science*(2*e*).

• My Own View

"All knowledge and skills of gaining and communicating insights from complex data through digital techniques. It is a blend of principles, algorithms, and systems to extract knowledge and insights from structured and unstructured data."

- Main contents of Data Science:
  - Data Collection: surveys, sensors, web scraping and OCR etc.
  - Data Wrangling: cleaning, transforming, merging, filtering, aggregating, and summarizing.
  - Data Analysis: descriptive statistics, causal inference, and predictive analytics.
  - Data Visualization: graphs, charts, and maps.
  - Data Communication: reports, dashboards, and presentations
- As a student that analyzes data, most of it you do already. However...

## What is Artificial Intelligence (AI)?

- AI, thus the *Artificial Intelligence*, is a field of computer science that aims to create machines that can perform tasks that typically require human intelligence.
- Many Areas:
  - Machine Learning
  - Natural Language Processing (NLP)
  - Computer Vision
  - Robotics
  - Expert Systems
- The most influential breakthrough in AI recently is in the space of **Generative AI** models or the **Large Language Models(LLM)** 
  - which are designed to create text or other forms of media based on patterns and examples they have been trained on such as ChatGPT and many others.
- It is dramatically changing the lives of human beings, and so is how we do research.

### Introduction to LLMs

- Large Language Models (LLMs) are one type of AI systems that generate human-like text.
- At their core, LLMs:
  - **Predict text** based on previous input.



### Introduction to LLMs

- Similar to other predictive models based on numerical data, LLM accuracy depends on:
  - Training data quality and quantity
  - Model complexity (measured by number of parameters)
- Training involves massive text datasets and complex neural networks.
  - For example, GPT-4 training utilized 45TB of data, over 300 trillion tokens, and involved 175 billion parameters.
  - Such training typically requires months and costs millions of dollars.
- Since LLMs are just *prediction model* of texts based on the input text, they **lack the ability to understand** the meaning of the text.

### Introduction to LLMs: Capabilities

• Though it does not really understand what you want, it can **generate** text that is **very similar** to what you want because it has been trained on a large amount of text data.

#### • Key Capabilities:

- Very knowledgeable and language-based.
- Processing diverse text-based tasks.
- Handling Math and Code.

#### • Extended Capabilities:

- Extend to multimodal tasks.
- Extend to lastest knowledge and information.
- Reasoning and inference capabilities.
- Extend to control applications, even devices.
- In summary, LLMs demonstrate remarkable **versatility** and **generality**, and most importantly, their capabilities continue to **evolve** at a rapid pace.

## Why Economics need Data Science and AI

• Social sciences(*firstly started by Economics*) have experienced two methodological revolutions over the past few decades.

#### • No.1: Credibility Revolution

- A movement that emphasizes the goal of obtaining secure causal inferences in social sciences.(Angrist and Pischke, 2010)
- The revolution started from around the 1990s, pioneering in economics, then spread over to other empirical social sciences such as sociology, political science, education, public policy, etc., which has entirely changed empirical social science and business research.

# The Sveriges Riksbank Prize in Economic Sciences in Memory of Alfred Nobel 2021



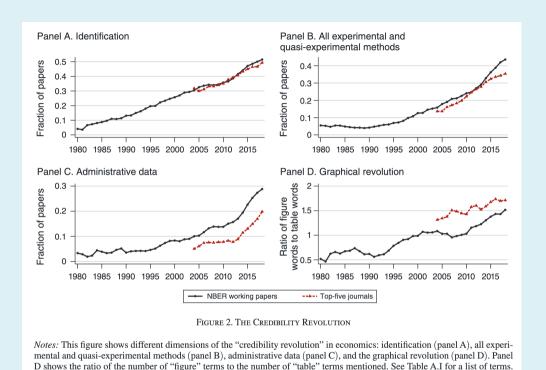
© Nobel Prize Outreach. Pho Paul Kennedy David Card Prize share: 1/2



© Nobel Prize Outreach. Pho Risdon Photography Joshua D. Angrist Prize share: 1/4



© Nobel Prize Outreach. Photo Paul Kennedy Guido W. Imbens Prize share: 1/4



Panel A. Difference-in-differences Panel B. Regression discontinuity 0.08 0.06 ction of I ⊆ 0.04 0.02 1990 1995 2000 2005 2010 2015 1990 1995 2000 2005 2010 2015 1980 1985 Panel C. Event study Panel D. Bunching ω 0.04 80.03 0.06 ₺ 0.04 ₲ 0.02 1980 1985 1990 1995 2000 2005 2010 2015 1980 1985 1990 1995 2000 2005 2010 2015 NBER working papers ---- Top-five journals FIGURE 4. QUASI-EXPERIMENTAL METHODS Notes: This figure shows the fraction of papers referring to each type of quasi-experimental approach. See Table A.I for a list

#### Key words for CR

The series show five-year moving averages.

#### Quasi-experimental methods

of terms. The series show five-year moving averages.

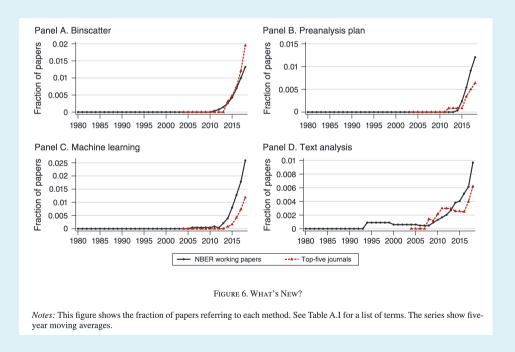
• Currie, J., Kleven, H., & Zwiers, E. (2020). Technology and Big Data Are Changing Economics: Mining Text to Track Methods. AEA Papers and Proceedings, 110, 42–48

#### • No.2: Big Data Revolution

 A movement that emphasizes that how our increasing ability to produce, collect, store and analyze vast amounts of data is going to transform our understanding of the human affairs. (Schonberger, 2013)



 Data sources and types are changing, which makes new methods to obtain, process, analyze and visualize data necessary.



Currie, J et al(2020)

- Now we are facing the third revolution in social science: No.3: AI Revolution.
- It is dramatically changing the way of obtaining, processing, analyzing and visualizing information and knowledge. So it is also changing the way of doing research.

	Usefulness
Brainstorming	•
Feedback	lacktriangle
Providing counterarguments	lacktriangle
Synthesizing text	•
Editing text	•
Evaluating text	•
Generating catchy titles & headlines	•
Generating tweets to promote a paper	•
Summarizing Text	•
Literature Research	0
Formatting References	•
Translating Text	•
Explaining Concepts	lacktriangle
y subjective rating of LLM capabilities as of Se	ptember 2023:
inconsistent and require significant human over	rsight
but will likely save you time	
ing this into your workflow will save you time	
	Feedback Providing counterarguments Synthesizing text Editing text Evaluating text Generating catchy titles & headlines Generating tweets to promote a paper Summarizing Text Literature Research Formatting References Translating Text Explaining Concepts y subjective rating of LLM capabilities as of Se inconsistent and require significant human over but will likely save you time

Category	Task	Usefulness
Coding	Writing code	•
	Explaining code	•
	Translating code	•
	Debugging code	•
Data Analysis	Creating figures	•
	Extracting data from text	•
	Reformatting data	•
	Classifying and scoring text	•
	Extracting sentiment	•
	Simulating human subjects	•
Math	Setting up models	•
	Deriving equations	$\circ$
	Explaining models	•
The third column repor	rts my subjective rating of LLM capabilities as	of September 2023:
O: experimental; result	s are inconsistent and require significant huma	n oversight
$\mathbb{O}$ : useful; requires ove	rsight but will likely save you time	
•: highly useful; incor	porating this into your workflow will save you	time

• Korinek, A. (2023). Generative AI for Economic Research: Use Cases and Implications for Economists.

## Why Economics needs Data Science and AI

- Economics and Data Science have many in common
  - require strong analytical skills
  - rely heavily on data and statistical analysis
  - need domain knowledge to interpret results

- The main difference between Economics and Data Science
- 1. Data collection and processing
  - Econ: more structured and cleaned
  - DS: more unstructured and messy
- 2. methods and tools
  - Econ: more theoretical and causal inference
  - DS: more practical and prediction

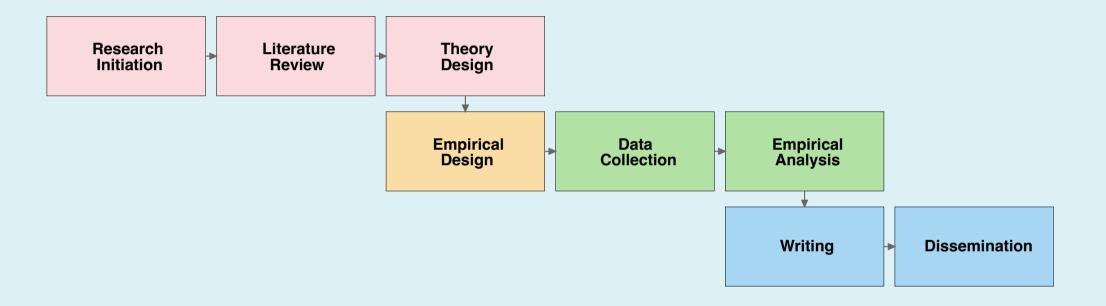
## Why Economics needs Data Science and AI

- Data Science can be seen as a complement to Economics, providing new tools and methods to analyze data and extract insights that may not be possible with traditional econometric methods alone.
- Economics theory and Econometrics methods can also provide a foundation for Data Science, helping to guide the analysis and interpretation of data.
- The combination of Economics and Data Science can lead to more robust and comprehensive analyses that can provide valuable insights into complex economic and social issues.

## Why Economics needs Data Science and AI

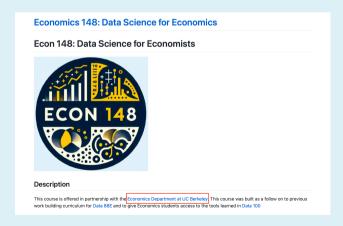
- AI revolution especially the **Generative AI** or **Large Language Models** are the result of the developing combination of Social Science, Data Science and Computer Science.
- It is dramatically changing the way we do research in Economics and Data Science
  - With the help with AI tools, we can finish the work more efficiently and accurately.
- Enhanced Research Efficiency: Accelerates literature reviews, theoretical modeling, data analysis, and academic writing.
- Expanded Research Perspectives and Scope: Facilitates interdisciplinary knowledge integration, identifies new research directions, and produces more comprehensive studies.
- Strengthened Methodological Tools: Provides novel approaches for research design, data collection, and methodological implementation.
- Advanced Academic Dissemination: Enables rapid dissemination of scholarly work and enhances research impact through AI-assisted communication.

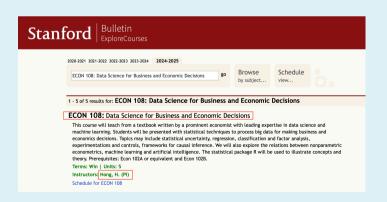
## LLMs in each step of the research workflow



## **Economics and Data Science Everywhere**

- Many double-majored programs in Economics and Data Science on both undergraduate and graduate levels have been established in many universities worldwide during the past 5 years
  - "The most popular double major in UC Berkeley"
- Many programs in Economics at top universities also provide the courses in Data Science and AI for their students.
  - MIT, Harvard, Stanford, Chicago, and UC Berkeley etc.





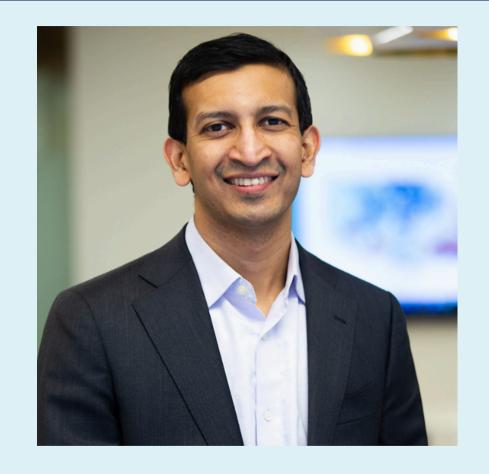
@Standford

@UCBerkeley

## The Most Popular undergraduate course

#### Top 3 in Harvard University

- Economics 10b: "Principles of Economics"
- Life Sciences 1b: "An Integrated Introduction to the Life Sciences"
- Economics 1152: "Using Big Data to Solve Economic and Social Problems"



Raj Chetty

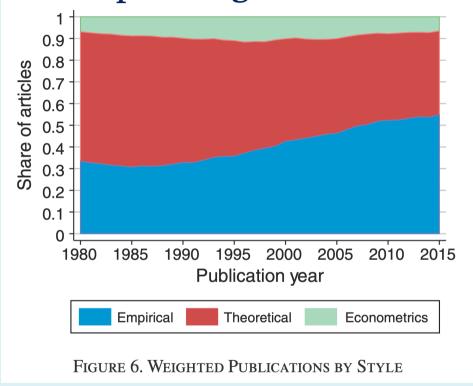
### Who'd better take the course?

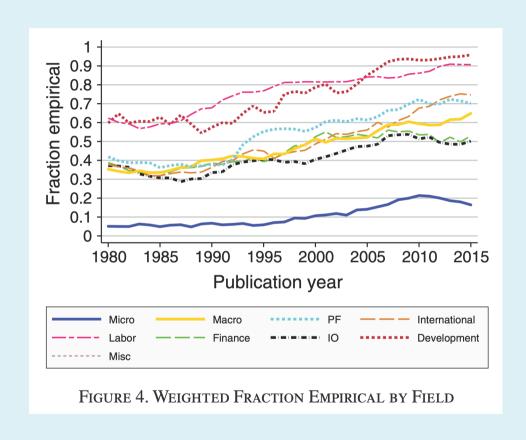
### Who' d better take the course?

#### The Purpose of the course

- Introduce you the foundational concepts of data science and artificial intelligence, emphasizing their practical applications in economics and social sciences.
- Unlike traditional econometrics courses, the focus is to learn some new tools and methods that can develop your ability to work with **non-traditional** economic data.
- The course is designed to be accessible to students with a wide range of backgrounds, including those with little or no prior experience in data science or economics.
- Help yourself enjoy to learn some new ideas in an empirical social scientist's mindset.

#### For those pursuing an academic career:





Angrist et al(2017)

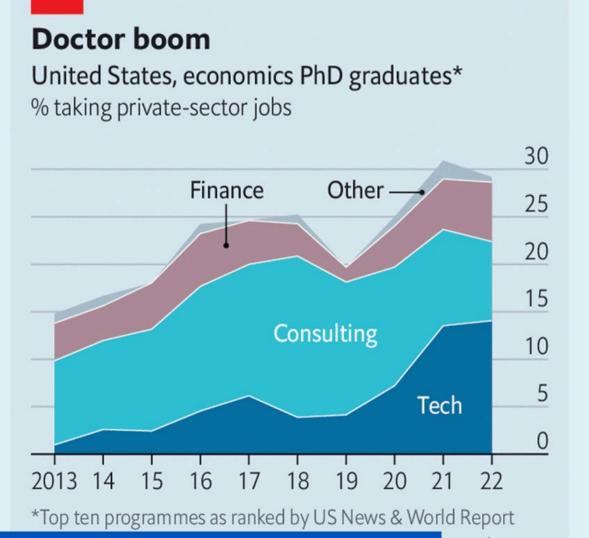
Angrist et al(2017)

• The proportion of empirical studies in economics is **increasing more and more**.

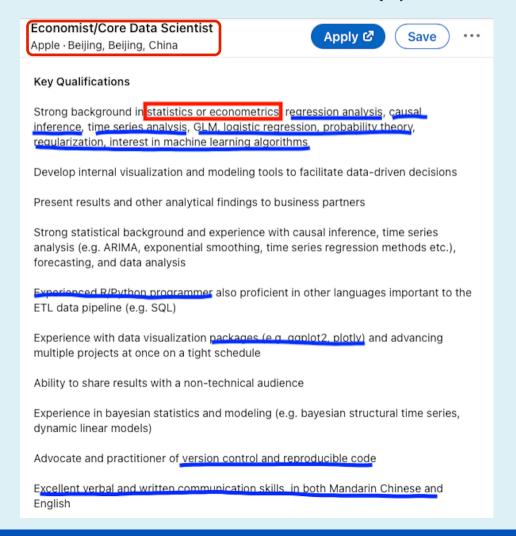
#### Who want to enter the industry job market

- Who want to work in industry: mastering econometrics may help you get a good job!
- A lot of internet giants even hire economists to lead their special R&D department. Such as
  - Google, Microsoft, eBay, Baidu, Alibaba, Tencent, Tiktok
- Data Scientist is the hottest job in consulting, business areas as well as financial industry right now.

The Change of EcoPhD Job Markets in the US by The Economist



### Who want to enter industry job market: Apple Job Wanted



#### Description

- Work with various teams to understand business problems and provide business solutions
- · Build models to causal impact of new programs release across different scenarios
- · Develop internal visualization and modeling to facilitate data-driven decisions
- · Present results and other analytical findings to business partners

#### **Education & Experience**

- PhD in Economics or related fields
- M.S. in related field with 5+ years experience applying econometric models to business problems.

#### Who want to enter industry job market: ByteDance Job Wanted

#### 国际电商-经济学家/数据科学家

上海 | 正式 | 产品 - 数据分析 | 职位 ID: A117677

#### 职位描述

我们欢迎有创造力、探索精神、且具备基本经济学、统计学素养的人才加入,和业务方共创并推动项目的上线落地。我们的 合作业务方包括推荐算法、产品、运营、资源管理等。

#### 主要职责:

把商业问题转化为可解的模型问题。通过经济学视角的思考和科学的方法(因果推断、AB实验、求解理论模型、预测等)来推动搜推策略、产品功能、资源分配等相关决策:

- 1、因果性的衡量各类策略、政策的效果,衡量长期影响,并形成系统性的方法论;
- 2、对数据现象现象进行归因,对用户、商家的决策链路做深入探索,总结洞察和建议,帮助各决策方建立认知;
- 3、优化各类资源分配(流量、营销补贴);
- 4、优化国际电商的生态环境,包括但不限于经营环境,用户体验、内容生态、供给生态、持续助力商达成长和用户增长。

#### 职位要求

- 1、经济学、统计学、运筹学、金融学、或者其他的相关的量化学科背景;
- 2、掌握 R 或Python等至少一项数据分析必备的编程语言,以及基础SQL能力;
- 3、有一定的解决商业问题、构建可落地的系统性解决方案、复杂项目管理、协调多方决策的经验;
- 4、良好的写作沟通能力;
- 5、以下领域的相关的科研、或者业界项目经历: reduced-form 因果推断、预测、causal ML、劳动经济学、健康经济学、教育经济学、行为经济学、金融经济学、产业组织学。

#### 商业化数据科学家-因果推断

上海 | 正式 | 产品 - 数据分析 | 职位 ID: A56405

#### 职位描述

- 1、通过积累日常使用经验、阅读相关学术论文和公开资料等,沉淀并向数科团队输出对因果推断方法论的深入理解和使用 经验、澄清常见的使用误区、提供标准应用流程指南、以保障方法在团队内应用的科学性、提高使用效率;
- 2、关注具有方法论共性或场景共性的相似业务问题,主导专项探索,与其他业务方向数科同学紧密配合,从宏观视角优化 资源分配效率或策略,优化产品策略,或针对相似问题抽象可复用的、普适的分析框架或解决方案,提升团队分析、决策效 变。
- 3、对宏观战略问题进行拆解、定义,通过数据描述、可视化、挖掘、统计建模等方法,提炼有效的数据洞察和产品战略建议、指导科学的决策与迭代。

#### 职位要求

- 1、本科以上学历,统计学、数学、计量经济学、数据科学、计算机等量化分析相关专业优先,硕士、博士优先;
- 2、具备扎实的统计学/计量经济学/机器学习/因果推断等数据科学理论基础及应用经验;精通SQL,熟练掌握Python/R中的一种,可进行数据清洗、可视化和分析;
- 3、具备快速学习能力,能够快速理解产品逻辑,并具备较强的逻辑思维能力,在较大不确定性的问题中可以构建分析框架、将数据转化为有效的商业洞察;
- 4、能够主动、独立思考的同时,具备良好的团队协作能力与责任心,善于与其他协作团队沟通,有主人翁意识;
- 5、具备强烈的好奇心与自我驱动力,乐于接受挑战,追求极致和创新,富有使命感。

投速

#### Who find a job in public sectors(shang'an movement)

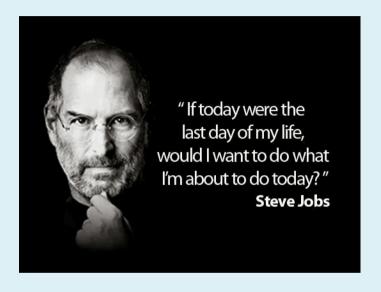


Source:SCMP

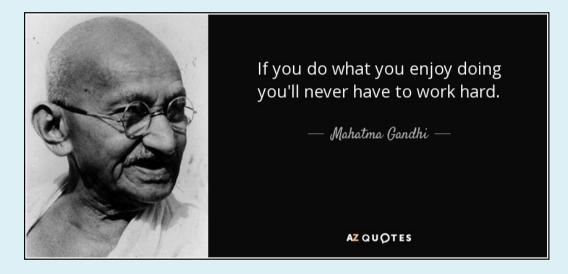
- To be honest with you, the course does not help you succeed in the examination.
- However, in the long run, it provides valuable knowledge and skills that offers a broader vision and abilities.
- It ultimately benefits your career, as well as our people, our country, and the world.

## Whoever and Whatever you want to be

#### Whoever you would like to be or whatever you want



• Every choice you make has an opportunity cost,try your best to make a wise one.



• Enjoy doing something seriously and cultivate a special quality for yourself!

#### Hard and Soft Skills

#### You COULD learn or improve several important skills during your college.

Hard Skills

Soft Skills

• Language

Critical Thinking

• Computer

• Teamwork

- Presentation and Writing
- Fortunately, you could learn/practice almost all above skills in our class.

## Wrap up

- In essence, DSAI is an cutting-edge and intriguing yet challenging course.
  - Please consider carefully before enrolling
  - Once committed, please work hard on it!
  - And remember, enjoy the process of working hard!

# **Course Logistics**

### About Me, our TA, and the Course

- My name is Zhaopeng Qu(曲兆鹏)
  - Associate Professor, Institute of Population Studies, Business School.
  - Research Fields: Labor Economics and Applied Econometrics
  - Email: qu@nju.edu.cn
- 南京大学"人工智能"建设系列课程:

  - 。 专业选修课 for 2nd/3rd year undergraduate students in Economics or other social sciences.
- 2024年秋季**第一次开课!**,2025年秋季**第二次开课!**, so everything is still fresh.
- Course Website: (https://byelenin.github.io/DSAI\_2025/)

- Our TA: Xiaotian Yu(虞晓天)
  - 2nd year PhD Student in Economics, Nanjing University.
  - Proficient in R and Python.
  - Will assist you in the lab sessions.
  - WechatID: yxtyxtii

## Prerequisite and Procedures

- Although there is no formal prerequisite for the course, I **recommend** that you'd better take at least one course in **Statistics**, **Data Analysis** or **Econometrics**.
- And I assume that you should be comfortable to dealing with data and coding experiences by using Python/R/Stata.
- The First Part: Lectures by the instructor
  - Introduce the underlying theoretical concepts briefly and focus on applications heavily.
  - Provide some specific examples in classical papers in the topics.
  - Normally, everyone who take the course is required attending the class.

- The Second Part: Coding practices in and after the class
  - Hands-on coding sessions to apply theoretical concepts.
  - Use of real-world datasets for practical experience.
  - Encourage collaboration and discussion among peers.

#### **Course Overview**

- Basic Tools
  - Github, IDE, AI tools, and more
- Introduction to R and Python for Data Analysis
  - Data wrangling, visualization, and analysis
- Topic 1: Spatial Data Analysis
  - Mapping, geocoding, and basic spatial analysis
  - NASA: Nightlight Data
- Topic 2: Web Scraping
  - Scraping websites, and some web index data

- Topic 3: OCR and Image Analysis
  - Extracting text from images, analyzing image data
  - Image recognition and classification
- Topic 4: Text Analysis
  - Text recognition and classification
  - Sentiment analysis
- Topic 5: Introduction to Big Data
  - Distributed computing frameworks (Hadoop, Spark)
  - Memory-efficient data processing techniques

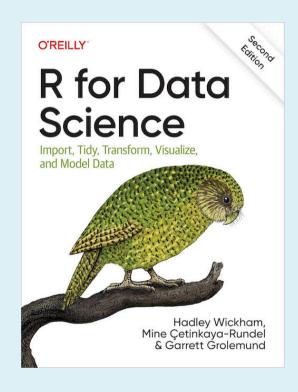
a Cloud computing solutions for economists

### Evaluation

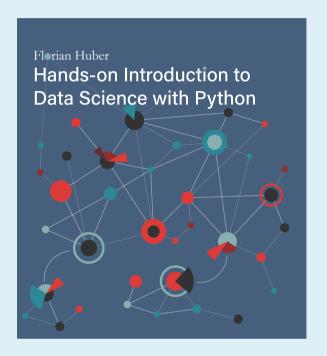
- The final grade will be based on the following components:
  - Class Participation(10%)
  - Proposals Presentation(40%)
  - Team Project Presentation(30%)
  - Final Report(20%)
- About Team Projects and Proposals
  - Students are required to form a team of 2-3 members to work on a research proposal.
  - Midterm Presentations: a presentation of the proposal.
  - Final Presentations: a presentation of the preliminary results.
  - Final Report: a final report of the project.

### **Textbooks**

• There are NO required textbooks for the course. However, the following textbooks are useful for reference:



• R for Data Science by Hadley Wickham, Mine Cetinkaya-Rundel and Garrett Grolemund (2e)



• Hands-on Introduction to Data Science with Python by Florian Huber

## **Computing Tools**

• The main computing tools used in the course are Python/R, optionally.

RPythonProConCon

- If your coding experience is limited to Stata, there's no need to worry.
  - I am currently learning Python and R as a beginner, and we can learn together.
  - Additionally, with the assistance of AI tools, anyone can become a proficient coder.

## Promise and Expectation

#### What I promise to offer you

- Prepare lectures as well as possible.
- Provide timely feedback on your projects.
- Provide additional resources and references.

- A good score?
  - It depends on you.

#### What I expect to you

- Class participation with a little bit aggressive attitude.
  - More questions, more scores!
- Self-motivated learning by doing.
  - More practices, more scores!

### Two Iron Rules



• Don't ever cheat on your assignments!



• Don't ever snitch your teachers to help political repression!

#### Welcome contact me



## Homework

### Homework

- Find a way to access a **stable and reliable internet** connection.
- Sign up to GitHub as a student and install Git on your computer.
- Then send your Github ID to out TA.