# Machine Learning for Economists
## Class 4: **Real** Financial Data

葛雷

中国人民大学经济院

2025 年 3 月 13 日

Real Financial Data

CSMAR & WRDS

Kaggle

Selenium

Homework3

# Real Financial Data

CSMAR & WRDS

Kaggle

Selenium

Homework3

# Real Financial Data

- What is the most important elements for Machine Learning? **Data**

- What makes the ML in finance unique? ( we financial data)

- Why real data?

## First look a the fake data

- sklearn.datasets is a good source for TOY data

- Good source for practice

- Only issue is that fake data is fake

- Lets check out why (Please follow to blank Ipynb)

Real Financial Data

# CSMAR & WRDS

Kaggle

Selenium

Homework3

# CSMAR

CSMAR, short for China Stock Market & Accounting Research
Database, is a comprehensive research-oriented database focusing
on China Finance and Economy. CSMAR was developed by
Shenzhen CSMAR Data Technology Co., Ltd based on academic
research needs, meeting with the international professional
standards while adapting to China's features.

# CSMAR

- professional level financial data for stock & company study

- used by both financial companies and financial researchers

# CSMAR: easy to use

- Easy to use especially for Python users

- We can use both UI and API (what is UA and API?)

- its check it with me step by step and login from lib

## WRDS

# USA counterpart of CSMAR [1]

---

[1]CSMAR followed WRDS's business model

# Kaggle

- Kaggle, a subsidiary of Google LLC

- Heavely platform for Quant Research (us)

- Codes, data, competition and more

- Let check it out! (Kaggle)

# Kaggle

- Kaggle is most important data source for now

- You can search and find your interested research topics

- Let check it out! (Kaggle)

Real Financial Data

CSMAR & WRDS

Kaggle

Selenium

Homework3

## Data from the internet

1. Internet has valuable data for the financial predictions

2. Internet data low quality? No

3. Selenium is a powerful and popular tool

# But how to use?

- I will guide you to study this package

- but next time you should know how to learn any package by yourself

# But how to use?

- Template + Documentation + CHATGPT + BING

- Template (from search bing and from CSDN, StackOverFlow, CHATGPT)

- 
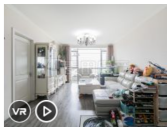  Something unknown → search Bing + Documentation + ChatGPT

# Homework3-1: Data Mining

- Housing Price Data from https://esf.fang.com/

- Housing Rent Data from https://zu.fang.com/

- Data needed: listed below

Real Financial Data
○○○

CSMAR & WRDS
○○○○○

Kaggle
○○○

Selenium
○○○○

Homework3
○○●○○○○

# Homework3-2: Data Mining

# Homework3-2: Data Mining

- Team 1 北京-海淀 I: 苏州桥、万柳、北太平庄、世纪城
- Team 2 北京-海淀 II: 西三旗、清河、西二旗、上地
- Team 3 河北-廊坊 + 北京-通州: 大厂、燕郊、马驹桥、亦庄
- Team 4 北京-昌平: 沙河、霍营、回龙观、天通苑
- Team 5 天津: 中新生态城 (滨海新区)、武清、劝业场 (和平)、八里台 (南开)
- Team 6 重庆-渝北: (Please choose blocks with both price and rental data)

- Each person only in charge of **one block** and only get first 20 pages if too many for you

# Homework3-3: Data Research

- Collect Data from your teammates and merge the data (please feedback to TA if someone no response, so we can help both team and other student)

- Data description of your data and whether data has outliers

- Then get housing price per m2 and housing rent per m2 ($price/m2$ and $rent/m2$) for each block

- Calculate median price to rent ratio for each block

- Figure A: Bar Plot the median price to rent ratio for each block (The global fair value should around 200)

## Homework3-4: Data Science

- Model 1
$price/m2_i = \beta_0 m2_i + \beta_2 location_i + \beta_2 m2_i \times location_i + \epsilon_i$

- Model2
$rent/m2_i = \beta_0 m2_i + \beta_2 location_i + \beta_3 m2_i \times location_i + \epsilon_i$

- Use model 1 and model 2 to predict price and rent for the $m2 = 50$, $m2 = 100$

- Figure B and C: Bar Plot the price to rent ratio for each block for the $m2 = 50$, $m2 = 100$

- Submission: only Ipynb codes to your personal folder (NO DATA PLEASE, Git is for codes not for data)