API



Demonstration and Vote! This will be included as part of your in-class participation.

I have developed an APP. Let's try it.



In your opinion, how does it get the temperature data?

https://www.youtube.com/embed/s7wmiS2mSXY?enablejsapi=1

Let me show you the basics.

```
1 import requests
2 url = "http://apis.juhe.cn/simpleWeather/query"
3 params = {
4         "city": "北京",
5         "key": "your API key",
6      }
7 response = requests.get(url, params = params)
8 data = response.json()
9 temperature = data['result']['realtime']['temperature']
10 print(temperature)
```

You need to apply for your API key (password) I am not sharing my key with you because a free key can only be used 50 times a day...

How to get the key?

Sign up an account here and activate it. Sign up for an API key.

ل ۇ	我的聚合	O	我的API
8	数据中心	^	
	我的API		+ 申请新数据
	我的数据块		
	我的收藏		⑦ 天气预报 ⅠD:73 免费 开通会员
÷	安全中心	~	AppKey: 38b0441a463012bee90bda232e0abedb 🚾 门 限每天50次请求,开通「黑钻会员」最高每天100000次数
ሐ	账号管理	~	📶 统计 🕟 测试 \Xi 文档 🕢 🖲 置顶
Ð	工单管理		
Ô	费用管理	~	
\$	应用中心		

Exercise

Apply for your own API key and develop your own weather APP!

```
1 import tkinter as tk
2 import requests
3 def get temperature(city):
       url = "http://apis.juhe.cn/simpleWeather/query"
       params = { "city": city, "key": "your key"}
       response = requests.get(url, params=params)
       data = response.json()
       temperature = data['result']['realtime']['temperature']
       spaces to add = 3 - len(temperature)
       output string = " " * spaces to add + temperature
10
11
       return output string
12 def update display(city):
       temperature = get temperature(city)
13
       display label.config(text="{}当前温度为: {}摄氏度".format(city, temperature))
14
15 root = tk.Tk()
16 root.title("City Temperature App")
17 display label = tk.Label(root, text=" 请选择你所在的城市 ", font=("Arial", 16))
18 display label.grid(row=0, column=0, columnspan=3, pady=10)
19 cities = ["北京", "上海", "香港"]
20 for index, city in enumerate(cities):
       city button = tk.Button(root, text=city, command=lambda c=city: update display(c))
21
       city button.grid(row=1, column=index, padx=10, pady=10, sticky="ew")
22
23 for i in range(3):
       root.grid columnconfigure(i, weight=1)
24
25
26 root.grid rowconfigure(1, weight=1)
27 root.mainloop()
```

Let's try the following API for Chinese/English translation

Translation

```
1 import requests
2 url = "https://api.oick.cn/api/fanyi"
3 params = {
4     "text": "欢迎来到香港大学",
5     }
6 response = requests.get(url, params = params)
7 data = response.json()
8 result = data['data']['result']
9 print(result)
```

Exercise: Design an APP so that users can input Chinese/English and the APP translates it into the other language

My APP created by GPT

```
1 import tkinter as tk
2 import requests
3 def translate(text):
       url = "https://api.oick.cn/api/fanyi"
       params = { "text": text }
       response = requests.get(url, params = params)
       data = response.json()
       result = data['data']['result']
       return result
10
11 def translate text():
       chinese text = input text.get("1.0", "end-1c") # Get text from input text box
12
       english text = translate(chinese text)
13
       input text.delete("1.0", "end") # Clear the input text
14
15
       input text.insert("1.0", english text) # Display translated text in the same text box
16
17 \text{ root} = \text{tk.Tk()}
18 root.title("Chinese to English Translator")
19 input text = tk.Text(root, height=10, width=50)
20 input text.pack(pady=10)
21 translate button = tk.Button(root, text="Translate", command=translate text)
22 translate button.pack(pady=10)
23 root.mainloop()
```

AviationStack is a platform which provides an API to track the status of flights. Its website is here: https://aviationstack.com/ You can find its API documentation here.

Exercise

- Acquire a free API Key from the platform
- Develop an APP, the APP takes flight number as input from the user, and displays its fight status, departing airport and arriving airport on the screen.
- Try to do everything yourself (with help from GPT)!

This is my APP.

Flight Status Checker	_		×	
Enter Flight Number:				
MU770				
Check Status				
Status: active				
Departing City: Singapore Changi				
Arriving City: Shanghai				

This is "my" AI (API)

Enter your question:		×
Do you like HKU?		
Ask		
As an AI, I don't have personal feelings preferences, so I don't "like" or "dislike anything, including institutions like the University of Hong Kong (HKU). However, I provide information about HKU! It is one most prestigious universities in Asia, kn its strong academic programs, research contributions, and vibrant campus life. In have specific questions about HKU, feel f	or e" of the own for f you ree to	^

Exercise

- Acquire an API Key from DeepSeek (https://platform.deepseek.com/)
- This is not free, but very cheap. You need to top up \$2 or RMB 10 to use the API.
- Then, develop an APP which takes a question from user, sends the question to the API, and displays the answer from API to the user.

Top up



API keys

Your API keys are listed below. The API key is only visible and can be copied once at creation. Save it securely. Do not share your API key with others, or expose it in the browser or other client-side code. In order to protect the security of your account, DeepSeek may also automatically disable any API key that we've found has leaked publicly. Usage of API keys created before April 25, 2024, was not tracked.

Name	Кеу	Created	Last used		
ximarketing	sk-16926*********************2e3e	2025-02-26	2025-02-26	0	Ŵ

If you need help on how to use the API, check the document here: https://api-docs.deepseek.com/

What are the main topics in modern economics?

Scraping made Simple

Scraping

Consider a simple URL: https://www.aeaweb.org/issues/785 American Economic Review is an academic journal in economics (we don't care anyway). Let's try to write a

scraper to download the article titles from the page:

Decisions under Risk Are Decisions under Complexity

by Ryan Oprea

(pp. 3789-3811)

Curbing Leakage in Public Programs: Evidence from India's Direct Benefit Transfer Policy

by Prabhat Barnwal

(pp. 3812-46)

Aiming for the Goal: Contribution Dynamics of Crowdfunding

by Joyee Deb, Aniko Öry, and Kevin R. Williams

(pp. 3847-76)

Scraping

Xi Li



I want to scrape article titles from https://www.aeaweb.org/issues/785 A title is like this: Decisions under Risk Are Decisions under Complexity Show me python code

23/02/2025, 10:08:28



GPT-40 (Tokens used 527)

To scrape article titles from the provided webpage, you can use Python with libraries such as requests and BeautifulSoup. Here's a sample code to get you started:

The code is a bit problematic...

```
1 import requests
 2 from bs4 import BeautifulSoup
 3 # URL of the page to scrape
 4 url = "https://www.aeaweb.org/issues/785"
 5 # Send a GET request to fetch the page content
 6 response = requests.get(url)
7 # Check if the request was successful (status code 200)
  if response.status code == 200:
       # Parse the page content with BeautifulSoup
       soup = BeautifulSoup(response.content, 'html.parser')
10
11
12
     # Find all article titles
13
       # Note: The exact HTML structure of the page needs to be inspected to
   find the correct tags and classes
       titles = soup.find all('h2', class ='title') # Example: Change these
14
   tags and classes based on actual HTML structure
15
16
      # Print the titles
17
      for title in titles:
           print(title.get text().strip())
18
19 else:
       print(f"Failed to retrieve the page. Status code:
20
   {response.status code}")
```

27

What's wrong with it?

```
1 import requests
 2 from bs4 import BeautifulSoup
 3 # URL of the page to scrape
 4 url = "https://www.aeaweb.org/issues/785"
 5 # Send a GET request to fetch the page content
 6 response = requests.get(url)
7 # Check if the request was successful (status code 200)
  if response.status code == 200:
       # Parse the page content with BeautifulSoup
       soup = BeautifulSoup(response.content, 'html.parser')
10
11
12
     # Find all article titles
13
       # Note: The exact HTML structure of the page needs to be inspected to
   find the correct tags and classes
14
       titles = soup.find all('h2', class ='title') # Example: Change these
   tags and classes based on actual HTML structure
15
16
      # Print the titles
17
      for title in titles:
18
           print(title.get text().strip())
19 else:
       print(f"Failed to retrieve the page. Status code:
20
   {response.status code}")
```

28

Open the webpage in your Chrome browser.

Click the upper right Chome setting button of your browser and you will be directed here.

	\$	D C	
ſ.	•		
c	New tab	(Ctrl+T
6	New window	0	trl+N
\$	New Incognito window	Ctrl+St	hift+N
8	Xi	Signed in	>
Θu	Passwords and autofill		>
5	History		>
Ŧ	Downloads		Ctrl+J
☆	Bookmarks and lists		>
88	Tab groups		>
Ð	Extensions		>
ŧ	Delete browsing data	Ctrl+Shit	lt+Del
e,	Zoom –	100% +	::
0	Print		Ctrl+P
Ð	Search with Google Lens		
₿ _E	Translate		
۵	Find and edit		>
C,	Cast, save and share		>
â	More tools		>

Choose "More tools"... Choose "Developer tools"...



Click the key button and you will get to "select an element in the page to inspect it".

Alternatively, use "Ctrl + Shift + C."

If needed, also click the button to switch between desktop and mobile version.



scroll

	<pre><article class="journal-article art_18081" id="10.1257/aer.20221227"></article></pre>
•••	<pre>Decisions under Risk Are</pre>
	Decisions under Complexity == \$0
	<pre><div class="article-item-authors"></div></pre>
	<div class="page-range">(pp. 3789-3811)</div>
	<pre><article class="journal-article art_18082" id="10.1257/aer.20161864"> </article></pre>
	<pre><article class="journal-article art_18083" id="10.1257/aer.20181851"> </article></pre>

The title information is here.

Upload your screenshot to ChatGPT!

Xi Li

I want to download article titles from https://www.aeaweb.org/issues/785 A title is like this: "Decisions under Risk Are Decisions under Complexity". Show me python code.

Capture.JPG JPG 44.23 KB

23/02/2025, 10:22:52



-40 (Tokens used 714)

Certainly! Below is a Python code snippet that uses the requests and beautifulsoup4 libraries to scrape the article titles from the specified webpage.

 $\overset{\circ}{\sim}$

The code I got

```
1 import requests
 2 from bs4 import BeautifulSoup
 4 # URL of the page to scrape
  url = 'https://www.aeaweb.org/issues/785'
   # Send a GET request to the page
  response = requests.get(url)
10 # Parse the page content with BeautifulSoup
11 soup = BeautifulSoup(response.content,
   'html.parser')
12
13 # Find all article titles
14 titles = soup.find all('h3', class ='title')
15
16 # Print each title text
17 for title in titles:
       print(title.get text(strip=True))
18
```

Scraping

Consider a simple URL: https://www.aeaweb.org/issues/785 Try to download the links of the articles now, with the help of GPT. Here is an example of the link:

https://www.aeaweb.org/articles?id=10.1257/aer.20221227

Scraping



My Code

```
1 import requests
```

```
2 from bs4 import BeautifulSoup
```

3 # URL of the page to scrape

```
4 url = 'https://www.aeaweb.org/issues/785'
```

```
5 # Send a GET request to the page
```

```
6 response = requests.get(url)
```

7 response.raise_for_status() # Check that the request was
successful

```
8 # Parse the page content with BeautifulSoup
```

```
9 soup = BeautifulSoup(response.content, 'html.parser')
```

```
10 # Find all article links
```

```
11 article_links = []
```

```
12 for a_tag in soup.find_all('a', href=True):
```

```
13 href = a_tag['href']
```

```
14 if href.startswith('/articles?id='):
```

```
15 full_url = 'https://www.aeaweb.org' + href
```

```
article_links.append(full_url)
```

```
17 # Print the article links
```

```
18 for link in article links:
```

```
19 print(link)
```

16

Exercise

- Go to https://www.aeaweb.org/issues/785 to download the links to each article.
- Go to each of the above links to download the abstract of the article.
- Print the all abstracts on the screen.

Exercise

First, I asked GPT to scrape the abstract for a given URL.

	XiLi
	I want to scrape the abstract from page https://www.aeaweb.org/articles?id=10.1257/aer.20221227 Abstract reads like this: https://www.aeaweb.org/articles?id=10.1257/aer.20221227 Show me python code
	Capture3.JPG JPG 94.53 KB
	23/02/2025, 18:11:24
GPT-40 (Tokens used 1419)	

Exercise I got the following code.

```
import requests
from bs4 import BeautifulSoup
```

```
# URL of the page to scrape
url = 'https://www.aeaweb.org/articles?id=10.1257/aer.20221227'
```

```
# Send a GET request to the page
response = requests.get(url)
response.raise_for_status() # Check that the request was successful
```

```
# Parse the page content with BeautifulSoup
soup = BeautifulSoup(response.content, 'html.parser')
```

```
# Find the abstract section
```

```
abstract_section = soup.find('section', class_='article-information abstract')
if abstract_section:
    # Extract and print the abstract text
    abstract_text = abstract_section.get_text(strip=True, separator=' ')
    print('Abstract:', abstract_text)
else:
```

```
print('Abstract section not found')
```

Exercise

I change it to a function "getabstract":

```
import requests
from bs4 import BeautifulSoup
def getabstract(url):
    response = requests.get(url)
    response.raise for status()
    soup = BeautifulSoup(response.content, 'html.parser')
    abstract section = soup.find('section', class ='article-information
abstract')
    if abstract section:
        # Extract and print the abstract text
        abstract text = abstract section.get text(strip=True, separator=' ')
        print('Abstract:', abstract text)
    else:
        print('Abstract section not found')
```

Exercise

Combine the previous codes, you got it!

```
1 import requests
 2 from bs4 import BeautifulSoup
   def getabstract(url):
       response = requests.get(url)
       response.raise for status()
       soup = BeautifulSoup(response.content, 'html.parser')
       abstract section = soup.find('section', class ='article-information abstract')
       if abstract section:
10
           # Extract and print the abstract text
11
           abstract text = abstract section.get text(strip=True, separator=' ')
           print('Abstract:', abstract text)
12
13
       else:
14
           print('Abstract section not found')
15
16
17 url = 'https://www.aeaweb.org/issues/785'
18 response = requests.get(url)
19 response.raise for status() # Check that the request was successful
20 soup = BeautifulSoup(response.content, 'html.parser')
21 article links = []
22 for a tag in soup.find all('a', href=True):
       href = a tag['href']
23
       if href.startswith('/articles?id='):
24
           full url = 'https://www.aeaweb.org' + href
25
26
           article links.append(full url)
27
28 for link in article links:
29
       getabstract(link)
```

We can do more!

- On page https://www.aeaweb.org/journals/aer/issues, we can scrape the URL of each issue.
- On the issue page, we can scrape the URL of each article.
- On the article page, we can scrape the abstract of each article (it took me two hours).
- Then, we can save all abstracts in one single file.
- After that, we can perform topic modeling!

We can do more!

I have already uploaded the txt file to the course website, which can be found here: http://ximarketing.github.io/data/AER_abstracts.txt

Please try the top words and LDA algorithm yourself and demonstrate the main topics in modern economics.

Exercise: Download the Photos of All HKU Business School Professors from https://www.hkubs.hku.hk/people/faculty? pg=1&staff_type=faculty&subject_area=all&track=professor iate®ion=hku

Downloading from a single page:

```
1 import requests
2 from bs4 import BeautifulSoup
3 import os
 4 url = "https://www.hkubs.hku.hk/people/faculty?
   pg=1&staff type=faculty&subject area=all&track=professoriate&region=hku"
5 response = requests.get(url)
 6 if response.status code == 200:
       soup = BeautifulSoup(response.text, 'html.parser')
       photo divs = soup.find all('div', class ='people-item')
10
       photo urls = []
12
       for div in photo divs:
           img tag = div.find('img')
13
           if img tag and 'src' in img tag.attrs:
14
               photo url = img tag['src']
15
16
               photo urls.append(photo url)
17
18
       os.makedirs('profile photos', exist ok=True)
19
       for idx, photo url in enumerate(photo urls):
20
21
           img response = requests.get(photo url)
22
           if img response.status code == 200:
23
               # Save the photo to the directory
24
               with open(f'profile photos/photo {idx + 1}.jpg', 'wb') as f:
25
                   f.write(img response.content)
26
               print(f'Downloaded: photo {idx + 1}.jpg')
27
           else:
28
29
               print(f'Failed to download photo from {photo url}')
30 else:
31
       print(f'Failed to retrieve the page. Status code: {response.status code}')
```

I further asked the following question...

I have 10 URLs each containing 16 photos, and I want to download all these photos. The URLs are <u>https://www.hkubs.hku.hk/people/faculty?</u> <u>pg=1&staff_type=faculty&subject_area=all&track=professoriate®ion=hku</u> to <u>https://www.hkubs.hku.hk/people/faculty?</u> <u>pg=10&staff_type=faculty&subject_area=all&track=professoriate®ion=hku</u>

I am done!





photo_2

photo_1



photo_11



photo_page1_5



photo_page1_15



photo_page1_16



photo_page1_7

photo_3

photo_13



photo_page2_1



photo_4

photo_14

photo_page1_8

photo_page2_2



photo_5



photo_15



photo_page1_9



photo_page2_3





photo_page2_4



photo_page1_10





photo_7

photo_page2_5









photo_8

photo_page1_12



photo_page2_6





photo_10

photo_page1_4



photo_page1_13



photo_page2_8

photo_page1_14























photo_6

















































The code I used

```
1 import requests
 2 from bs4 import BeautifulSoup
 3 import os
 4 base url = "https://www.hkubs.hku.hk/people/faculty?pg=
   {}&staff type=faculty&subject area=all&track=professoriate&region=hku"
 5 os.makedirs('profile photos', exist ok=True)
 6 for page in range(1, 11):
       url = base url.format(page)
       print(f"Fetching page: {url}")
       response = requests.get(url)
       if response.status code == 200:
10
           soup = BeautifulSoup(response.text, 'html.parser')
11
           photo divs = soup.find all('div', class ='people-item')
12
           for idx, div in enumerate(photo divs):
13
14
               img tag = div.find('img')
               if img tag and 'src' in img tag.attrs:
15
                   photo url = img tag['src']
16
                   img response = requests.get(photo url)
17
                   if img response.status code == 200:
18
19
                       photo filename = f'profile photos/photo page{page} {idx + 1}.jpg'
20
                       with open(photo filename, 'wb') as f:
21
                            f.write(img response.content)
                       print(f'Downloaded: {photo filename}')
22
23
                   else:
                       print(f'Failed to download photo from {photo url}')
24
       else:
25
           print(f'Failed to retrieve the page. Status code: {response.status code}')
26
```

Try to scrape article titles from the following page: https://pubsonline.informs.org/toc/mnsc/71/1

What else can we do with scraping?

What else can we do with scraping?

- Scraping a website with selenium.
- Web scraping with concurrency (parallel scraping).
- Changing IP addresses with VPN (using proxies/VPNs).

https://www.youtube.com/embed/cc26zFE8X1k?enablejsapi=1

Scraping with Selenium

You need to install selenium and download Chrome driver Moreover, add Chrome driver to your system path

```
1 from selenium import webdriver
2 from selenium.webdriver.common.by import By
3 url = "https://pubsonline.informs.org/toc/mnsc/71/1"
4 driver = webdriver.Chrome()
5 driver.get(url)
6
7 titles = driver.find_elements(By.CLASS_NAME, 'issue-item__title')
8 # Extract and print the text of each title
9 for title in titles:
10 article_title = title.find_element(By.TAG_NAME, 'a').text
11 print(article_title)
12 # Close the WebDriver
13 driver.implicitly_wait(5)
14 driver.quit()
```

Class Review

Basic Operations

Linear Regression Logistic Regression Multinomial Logit Model Artificial Neural Network Data Visualization Statistical Tests

Advanced Algorithms

Clustering (and Expectation-Maximization) Monte Carlo Algorithm Dynamic Programming Divide and Conquer

Unstructured Data

Text Data: Word Frequency, Topic Modeling, Summarization Image: Color, Brightness, Convolution, etc. You can also find APIs for text and image analysis.

API and Scraper

How to apply for and use an API. How to rely on AI to generate your own data scraper.