第3章

ChatGPT 与 Python 数据采集

数据采集与办公自动化有着密切的关系,两者的关系如下:

(1)数据源获取:办公自动化一般需要从不同的数据源中获取数据,如数据库、API、文件等。数据采集负责从这些数据源中收集数据,以供后续的自动化处理和分析使用。

(2) 实时数据更新:办公自动化可能需要基于实时数据进行操作和决策。数据采集可 以定期或实时地从数据源中获取最新的数据,确保自动化过程始终基于最新的数据进行 操作。

(3)数据预处理:数据采集包括对采集到的数据进行预处理和清洗的过程。办公自动 化可能需要处理缺失值、去除异常值和标准化数据等。数据采集可以在采集数据的同时进 行必要的数据预处理,以提供高质量的数据供办公自动化使用。

(4)数据整合和转换:在办公自动化过程中,可能需要整合来自不同数据源的数据,或 将数据转换为所需的格式。数据采集负责从不同数据源中收集数据,并进行整合和转换,以 满足办公自动化的需求。

(5)决策支持和报告生成:办公自动化一般需要基于数据进行决策支持和报告生成。 数据采集可以通过收集和整理相关数据,为自动化过程提供必要的数据基础,从而帮助生成 决策支持材料和报告。

综上所述,数据采集在办公自动化中起着重要的作用,它为自动化过程提供了数据的基础和支持。通过数据采集,办公自动化可以更准确、高效地处理和分析数据,帮助提高工作效率和决策质量。

3.1 ChatGPT 在数据采集过程中的作用

ChatGPT 在数据采集过程中的作用如下:

(1) 解析和处理数据采集的结果: ChatGPT 可以用于解析和处理数据采集的结果。用 户可以与 ChatGPT 交互,让它帮助解释和处理数据采集得到的结果,例如提取特定字段、 过滤数据、转换格式等。

(2) 提供数据采集过程中的帮助和指导: ChatGPT 可以作为一个智能助手,在数据采

集过程中提供帮助、建议和指导。用户可以向 ChatGPT 提出问题,寻求关于数据采集的最 佳实践、技巧或解决方案。

然而,实际上,ChatGPT 作为一个语言模型,它无法直接进行网络请求、数据抓取或处理。为了实现自动化的数据采集和处理,需要结合 ChatGPT 与其他适当的 Python 库和工具,例如 Web Scraping 库、API(应用程序编程接口)调用库、数据库访问库等。

这种结合可以帮助实现自动化的数据采集过程,通过编写 Python 代码和与 ChatGPT 的集成,实现数据的获取、处理和存储。

总结来说,ChatGPT 在数据采集过程中可以作为辅助工具和智能助手,但其本身并不 直接处理和执行数据采集任务。它可以提供对话交互和文本生成的能力,帮助用户在数据 采集过程中做出决策、提供反馈和解决方案。具体的自动化数据采集任务的实现还需要结 合其他的数据采集工具和技术。

3.2 数据采集概述

数据采集是指通过各种技术手段和方法,从不同的数据源中获取数据的过程。数据采 集在许多领域中都非常重要,包括科学研究、市场调查和业务决策等。随着互联网的普及和 大数据时代的到来,数据采集变得更加广泛和重要。

3.2.1 数据采集的重要性和挑战

数据采集的重要性在于它提供了信息和洞察力,可以用于支持决策、分析和创新。通过 采集各种类型的数据,我们可以发现模式、趋势和关联,从而得出有用的结论。数据采集还 可以帮助建立和改进模型、算法和预测,从而提高工作效率和结果准确性。

然而,数据采集也面临以下挑战。

(1)数据来源的多样性:数据可以来自各种不同的来源,如网页、数据库、传感器等。 不同的数据源可能使用不同的格式和协议,需要采集和整合这些数据变得更加复杂。

(2)数据量的庞大:随着数据产生速度的加快和存储成本的降低,我们能够处理的数据量越来越大。采集和处理大规模数据集需要具备高效的算法和技术。

(3)数据质量的保证:采集到的数据可能存在噪声、缺失或错误。确保数据的质量对 于后续的分析和应用非常重要,需要进行数据清洗和验证的工作。

3.2.2 数据采集的基本步骤

数据采集的基本步骤可以概括为以下几个阶段:

(1) 定义采集目标:明确你想要采集的数据类型、范围和目标。这有助于确定采集的 重点和方向。

(2)确定数据源:确定你需要采集的数据来源,如网页、API、数据库等。了解数据源的 结构、格式和访问方式。 (3)设计采集方案:根据数据源的特点和要求,设计采集方案。这包括选择合适的采 集工具、编写采集代码和设置采集参数等。

(4)执行数据采集:根据采集方案,执行数据采集操作。这可能涉及网络爬虫、API调用、数据库查询等技术手段。

(5)数据清洗和整合:采集到的数据可能需要进行清洗和整合,以确保数据的质量和 一致性。这包括处理缺失值、去除噪声、统一格式等操作。

(6)存储和管理数据:将采集到的数据进行存储和管理,以便后续的分析和应用。可 以使用数据库、数据湖、云存储等工具和技术来管理数据。

3.2.3 数据采集技术和工具

在数据采集过程中,可以使用一些常见的技术和工具。

(1) 网络爬虫: 通过自动化程序从网页中提取数据。

(2) API 调用: 使用 API 与远程服务器进行通信,获取数据。

(3)数据库查询:通过执行 SQL 查询语句从数据库中检索数据。

(4) 传感器和设备:使用传感器和设备采集实时数据,如温度、湿度、位置等。

(5)数据处理框架和库:使用 Python 中的数据处理框架和库(如 Pandas、NumPy、 BeautifulSoup 等)处理和分析采集到的数据。

这些技术和工具可以根据具体的需求和情况进行选择和应用。

3.3 网页数据爬取

从网页中爬取数据是数据采集的一种常见形式。在网页爬取数据时,可以使用 Python 编程语言结合相应的库和工具来自动化地从网页中提取数据。下面介绍一些常用的网页爬 取数据技术和方法。

3.3.1 网页中的静态和动态数据

裹挟在 HTML 代码中的数据并非唾手可得。大多数情况下,Web 前端与后台服务器进行通信时采用同步请求,即一次请求返回整个页面所有的 HTML 代码,这些裹挟在 HTML 中的数据就是所谓的"静态数据"。为了改善用户体验,Web 前端与后台服务器通 信也可以采用异步请求技术 AJAX^①,异步请求返回的数据就是所谓的"动态数据",异步请 求返回的数据一般是 JSON 或 XML 等结构化数据,Web 前端获得这些数据后,再通过 JavaScript 脚本程序动态地添加到 HTML 标签中。

① AJAX(Asynchronous JavaScript and XML)可以异步发送请求获取数据,请求过程中不用刷新页面,用户体验好,而且异步请求过程中,不是返回整个页面的 HTML 代码,只是返回少量的数据,这可以减少网络资源的占用,提高通信效率。



同步请求也可以有动态数据。就是一次请求返回所有 HTML 代码和数据,数据并不 是 HTML 放到标签中的,而是被隐藏起来,例如放到 hide 等隐藏字段中,或放到 JavaScript 脚本程序的变量中。然后通过 JavaScript 脚本程序动态地添加到 HTML 标签中。

图 3-1 所示的搜狐证券网页显示了某只股票的历史数据,其中图 3-1(a)所示的 HTML 内容都是静态数据,而动态数据则由 JavaScript 脚本程序动态地添加到 HTML 标签中,如 图 3-1(b)所示。

M Apple Inc. Common 5	* 4 BHRE:1139.79-0.40 *			e - a ×	Apple Inc. Common 5	× 4 8HF	B:1139.79 -0.4	× =						•	
- C O 不安全	q.stock.sohu.com/cn/600519/ishq.shtml			ŵ 🖸 i	← · C ◎ 不安全	q.stock.sohu.c	om/cn/60051	9/isbq.shtml	È.						* •
Senu.com RELize	◎ 按账第页-新闻-体育-S-	媒乐-r-射组-II-汽车-房产-	家第一女人-我想-董喜-Chirale	m-邮件-博客-加S-批说两切-建筑	Selfu.com RELES	F	按紧紧闭-	新闻-体育-\$-	编乐-v-新辑	- 〒- 汽车-	向市-家園	-女人-祝	频-频客-Chinake	m-邮件-请高-185-	- 我说的句- 建筑
投级证券 IIII stock sohu.com	10.6114.0mm 电电力中等机制用器电路中 在英一个级 指数	in . Hiti bi	ik nan	● 取入税税用 ● 逆力制作	搜狐 证券 stock sofu com	1070-14 0805 首 宛	単見た所 (个 取	≦81.981738388 87 %	198 	ŧŕ	板族		€ stan	(1) 加入的副兵	(点) 股为财政
代码/名称/带马	(152/929)	3 . нас.	101 (M)	业学和在8 iit	代码/名档/页写	TEAN			3	101.11		10	100	270	20 11
最近刘克隆 我的教徒教	<u> 新州茅台</u> 1139.79 单0.46% - (600519) 2020-04-03 15:00	5.21 S Ref 194	stalat o		最近30花板 我的自己来	<u>東州茅台</u> (600519)	1139.79	\$ 0.46% ·	-5.21		时行情	0.80	nam デカ太	米石七星	
	案时行情 成交明细 分约	未 历史计辑	北皮博教育 大家交易			实时行播	成文明建	95	18	历史行播	2	1.0.192	4/1元	,女义1/h	
	历史行情					历史石情						1			
	H. 2020-2-1	至 3030-04-03	◎ 接目 ◎ 接用 ● 接月	应 进			H 2	021-2-1	2	2000-04-03		10	日回使用日使用	重调	
2 BHF0					② 费州茅台	BM	开盘	0.2	ana a	#11 %	#1	1=	威交量(有)	滅交金数のう	NFB
						累计:	2020-02-05星	2020-04-03	86.99	8.208	900.1	1176	2167627	22485922.4	10.08
行情则表					行情能表	2020-04-05	1129.05	1109, 79	-6.23	-0,406	\$131.98	1147.96	27638	314942.19	0,238
the second se					实时行情体等	2020-04-02	1504.00	1145.00	42,00	7, 828	1103.88	1145.00	45205	50952.6.75	0.508
实时行情""									-8.00	-0.725	1105.00	1128.00	33205	370016.44	0,208
其町行情 ¹¹⁰ 成交明组					際の情報	2020-04-05	1137.90	1107 00	11.00						
或可行情 成交明细 分价表 历史经增					版文明组 分价表 ● 历史行情	2020-04-01 2020-03-31	1137.90	1103.00	29,00	2.646	1081.80	1115.00	47304	529317.06	0, 208
或可行情 中 成交明组 分价表 • 历史行情 文庫時前編					 成 2 件地 分 計表 6 安行情 2 序級教務 	2020-64-68 2020-63-38 2020-63-30	1137,90 1082,00 1090,25	1101.00 1111.00 1072.00	19, 00 -1, 50	3. 646 -0. 336	1081.80 1057.00	1115.00 1077.00	47104 30687	529317.06 327982.05	0,265
英町行簿 4 4 4 2 2 3 4 2 4 4 4 4 4 4 4 4 4 4 4 4					成之州坦 分供表 ▶ 历史行情 龙虎禄数据 大衆交易	2020-04-01 2020-03-31 2020-03-30 2020-03-27	1137.00 1082.00 1080.25 1085.00	1105.00 1111.00 1072.00 1075.50	-1.50 10.58	3.646 -0.336 0.996	1083.80 1057.00 1075.01	1115.00 1077,00 1092,00	47164 20087 20086	529317.40 327982.51 995175.94	0, 288 0, 248 0, 298
與时行情 ^{化44} 成交明组 分位委 • 历史行情 龙虎排教指 大奈交易					成交州班 分析表 ● 历史行情 龙虎特数据 大宗交易	2020-04-01 2020-03-38 2020-03-30 2020-03-27 2020-03-26	1137,00 1082,00 1085,00 1095,00 1073,33	1101.00 1111.00 1072.00 1075.50 1064.92	10,00 -1,50 10,59 -15,18	3, 848 -0, 338 0, 998 -1, 418	1083.80 1057.00 1075.01 1082.00	1115.00 1077,00 1092,00 1091,20	47904 30687 36202 00048	529317.00 327982.31 995175.94 525375.44	0, 268 0, 248 0, 298 0, 268
與时行情 成交明细 分价表 ● 历史行情 龙虎排放缩 大家交易 翻阅赞机					 ★ 2 / 1 / 1 / 1 / 1 / 1 / 1 / 1 / 1 / 1 /	2020-04-01 2020-03-38 2020-03-30 2020-03-27 2020-03-25 2020-03-25	1137, 00 1082, 00 1080, 25 1085, 00 1075, 33 1087, 98	1101.00 1111.00 1072.00 1075.50 1064.92 1080.10	-5, 50 -5, 50 10, 58 -15, 18 26, 10	3.846 -0.238 0.998 -1.428 2.228	1083.80 1057.00 1075.01 1062.00 1074.20	1115.00 1077.00 1092.00 1093.20 1093.36	47964 30087 36205 20049 27049	52903.7.00 327902.5s 399375.94 522375.44 403505.12	0, 208 0, 248 0, 298 0, 298 0, 298
实时行情 成交明组 分论表 • 历史行情 龙虎排数据 大 教文品 翻 阅货以 公司新闻					成义师地 分析表 下 历史行情 之虎梯数据 大和交易 載問要派 公明新闻	2020-04-01 2020-00-38 2020-00-30 2020-00-27 2020-00-25 2020-00-25 2020-00-25	1137, 80 1082, 00 1080, 25 1095, 00 1073, 33 1087, 98 1083, 00	1101.00 1111.00 1072.00 1075.50 1094.92 1080.10 1055.00	-1, 50 -1, 50 10, 58 -15, 18 24, 10 37, 00	3.646 -0.338 0.998 -1.418 2.228 3.538	1083.80 1057.00 1075.00 1062.00 1074.20 1035.00	1115.00 1077.00 1092.00 1091.20 1091.36 1088.88	47984 30687 36264 50049 37049 46655	529431.46 325982.52 399375.94 525375.44 461596.12 490335.44	0, 288 0, 248 0, 298 0, 298 0, 298 0, 379
與可行機 成交明細 分化表) 历史行情 龙虎导致器 大家交易 動陶賽紙 公司新闻 公司公告					成交相由 分析表 • 历史行情 龙虎特教授 大家交易 數何赞讯 公司新闻 公司公告	2020-04-05 2020-03-33 2020-03-30 2020-03-27 2020-03-25 2020-03-25 2020-03-25 2020-03-24 2020-03-25	1117, 00 1082, 00 1080, 25 1085, 00 1073, 33 1087, 98 1083, 00 1083, 00	1101.00 1111.00 1072.00 1075.50 1054.92 1080.10 1055.00 1019.00	-16.00 38.00 -1.50 10.59 -15.18 24.10 37.00 -16.28	3, 646 -0, 228 0, 998 -1, 418 2, 228 3, 638 -1, 578	1083.80 1057.00 1075.01 1082.00 1074.20 1035.00 991.52	1115.00 1077.00 1092.00 1093.20 1093.35 1088.00 1005.20	47984 30687 36566 30049 27049 27049 49655 49855 49855	52943.7.00 3227482.32 3995175.34 3525375.44 4615305.12 4900356.44 491,812.19	0, 288 0, 248 0, 298 0, 298 0, 298 0, 378 0, 378

(a) 静态数据

(b) 动态数据



3.3.2 使用 urllib 爬取静态数据

urllib 是 Python 标准库中的一个模块,提供了用于进行 HTTP 请求的基本功能。它包含了多个子模块,用于不同的请求任务和操作。

下面是 urllib 库中主要的子模块及其功能:

(1) urllib. request:用于发送 HTTP 请求和获取响应。它提供了一些函数,如 urlopen()用于打开 URL 并返回响应对象,urlretrieve()用于下载文件等。

(2) urllib. parse:用于解析 URL、拼接 URL 和处理 URL 编码。它包含了一些函数,如 urlparse()用于解析 URL 字符串,urljoin()用于拼接 URL,urlencode()用于将参数编码为 URL 查询字符串等。

(3) urllib. error: 定义了与 URL 请求相关的异常类。当在请求过程中发生错误时,可以捕获这些异常来进行适当的处理。

(4) urllib. robotparser:用于解析和分析 robots. txt 文件,该文件用于指示爬虫哪些页

面可以访问。

使用 urllib 库,你可以发送 HTTP 请求(GET、POST 等)、设置请求头、处理响应数据和错误等。它是 Python 标准库的一部分,因此不需要安装额外的依赖库。

然而,需要注意的是,urllib 库相对较低级,对于一些高级的功能,如处理异步请求、处理复杂的请求参数和响应处理,可能需要借助其他库,如 requests 或 aiohttp。

下面是一个简单的示例,演示了使用 urllib. request 发送 GET 请求并获取响应的过程:

import urllib.request

url = 'https://example.com'
发送 GET 请求并获取响应
response = urllib.request.urlopen(url)
读取响应内容
data = response.read()
关闭响应
response.close()
处理数据
...

3.3.3 示例 1: 爬取纳斯达克苹果股票数据

下面通过一个案例介绍如何使用 urllib 爬取静态网页数据,如图 3-2 所示是纳斯达克 苹果公司股票历史数据网页。

```
示例实现代码如下:
# 代码文件:chapter3/3.2.2.py
import os
import urllib. request
# url = 'https://www.nasdaq.com/symbol/aapl/historical#.UWdnJBDMhHk'
                                                                                  1
# 本地文件访问
url = "file:///" + os.path.abspath("./nasdag-Apple1.html").replace("\\", "/")
                                                                                  (2)
req = urllib.request.Request(url)
with urllib.request.urlopen(req) as response:
  data = response.read()
  html data = data.decode()
  print(html data)
示例运行后,输出结果如下:
<! doctype html >
< html lang = "en">
< head >
```

Apple Inc. Commo	n Ste 🗙 📃						θ - Ο
$- \rightarrow \mathbf{C}$ O https	://www.nasda	aq.com/sy	mbol/aapl/l	historical#.U	WdnJBDMhHk		6 \$
Nasdaq		E	nter symbo	l, name or k	eyword	Sear	ch Join NASDAQ or Log
	QUOTES V	MARKET	's 🔻 / NEW	s v / INVE	STING T ADV	ANCED INVESTIN	G V PERSONAL FINANCE V MY NASDAQ V
SYMBOL LIST VIEWS	Results for:	3 Month, F	rom 15-DEC	-2017 TO 15	-MAR-2018		Interactive Brokers Debit Mastercard®
InfoQuotes	Date	Open	High	Low	Close / Last	Volume	Debit Mastercard
TOCK DETAILS	16:00	178.50	180.24	178.0701	178.65	22,743,798	
Summary Quote	03/15/2018	178.5	180.24	178.0701	178.65	22,676,520	
Real-Time Quote	03/14/2018	180.32	180.52	177.81	178.44	29,334,630	
After Hours Quote	03/13/2018	182.59	183.5	179.24	179.97	31,464,170	
Pre-market Quote	03/12/2018	180.29	182.39	180.21	181.72	32,191,070	AutoractiveBrokers
option Chain	03/09/2018	177.96	180	177.39	179.98	32,130,360	
HARTS	03/08/2018	175.48	177.12	175.07	176.94	24,518,850	
lasic Chart	03/07/2018	174.94	175.85	174.27	175.03	31,686,450	5314 0012 3456 7890 Deter
nteractive Chart	03/06/2018	177.91	178.25	176.13	176.67	23,273,160	T. PETERFFY modeled
OMPANY NEWS	03/05/2018	175.21	177.74	174.52	176.82	27,825,140	Lise your securities as collatoral to
Company Headlines	03/02/2018	172.8	176.3	172.45	176.21	38,426,060	borrow at 2.66% APR* or less.
ress Releases	03/01/2018	178.54	179.775	172.66	175	48,706,170	1
larket Stream	02/28/2018	179.26	180.615	178.05	178.12	37,568,080	InteractiveBrokers
TOCK ANALYSIS	02/27/2018	179.1	180.48	178.16	178.39	38.885,510	OPEN AN]
nalyst Research	02/26/2018	176.35	179.39	176.21	178.97	37 353 670	L ACCOUNT J
uru Analysis	00/00/00 40	470.03	475.05	470.54	475 5	00 770 050	Member - NYSE, FINRA, SIPC. For more info see ibkr.com/int
Stook Report	02/23/2018	173.67	175.65	173.54	1/5.5	33,112,050	

图 3-2 纳斯达克苹果公司股票历史数据网页

```
< meta charset = "UTF - 8">
 < meta name = "Generator" content = "EditPlus">
 < meta name = "Author" content = "">
 < meta name = "Keywords" content = "">
 < meta name = "Description" content = "">
 < title > Document </title >
</head>
< body >
< div id = "quotes content left pnlAJAX">
 < thead class = "historical - data__table - headings">
    Date 
     Open 
     High 
     Low 
     Close/Last 
     Volume 
   </thead>
```

```
>10/04/2019
     225.64 
    227.49
    223.89
    227.01
    34,755,550 
  >10/03/2019
    218.43
     220.96 
   215.132 
    220.82 
    30,352,690 
  </div>
</body>
</html>
代码解释如下:
```

- •代码第①行指定 URL 网址。
- •代码第②行中指定本地文件地址 nasdaq-Apple1.html。



为什么要采用本地文件呢?这是因为我们爬取的网址经常容易改版,出于学习方便,笔者提供了本地文件,读者需要注意根据自己的实际情况将代码第②行改成自己的文件地址。

3.4 解析数据

数据爬取回来后,需要从 HTML 代码中分析出需要的数据,这个过程可以使用适当的 数据解析技术实现,例如使用正则表达式、BeautifulSoup 和 XPath 等进行 HTML 或 XML 解析,或使用 JSON 解析库处理 JSON 数据。笔者推荐使用 BeautifulSoup 库。当然也可以 利用 ChatGPT 辅助解析数据。

3.4.1 使用 BeautifulSoup 库

BeautifulSoup 是一套帮助程序设计师解析网页结构项目, BeautifulSoup 官网是 https://www.crummy.com/software/BeautifulSoup/。

要使用 BeautifulSoup 库首先需要安装 BeautifulSoup,可以通过如下 pip 指令进行 安装:

pip install beautifulsoup4

安装过程如图 3-3 所示。

🔜 命令提示符



图 3-3 BeautifulSoup 安装过程

下面我们介绍 BeautifulSoup 常用 API。

BeautifulSoup 中主要使用的对象是 BeautifulSoup 实例, BeautifulSoup 常用函数如下:

- find_all(tagname): 根据标签名返回所有符合条件的元素列表。
- find(tagname): 根据标签名返回符合条件的第一个元素。
- select(selector):通过 CSS 中选择器查找符合条件所有元素。
- get(key, default=None):获取标签属性值,key是标签属性名。

BeautifulSoup常用属性如下:

- title: 获得当前 HTML 页面的 title 属性值。
- text: 返回标签中的文本内容。

在网络爬虫抓取 HTML 代码时,开发人员需要知道数据裹挟在哪些 HTML 标签中, 要想找到这些数据,可以使用一些浏览器中的 Web 开发工具。笔者推荐使用 Chrome 或 Firefox 浏览器,因为它们都自带了 Web 开发工具箱。Chrome 浏览器可以通过菜单"更多 工具"→"开发者工具"打开,如图 3-4 所示。Firefox 浏览器可以通过菜单"Web 开发者"→ "切换工具箱"打开,如图 3-5 所示。还可通过快捷键打开它们,在 Windows 平台下两个浏 览器打开 Web 工具箱的快捷键都是 F12。

Apple Inc. Commo	on St. X								0 - 0	
← → C ● 安全	https://www	w.nasdaq.o	.om/symbol	/aapl/histor	ical#.UWdnJB	DMhHk			ជំ	• :
Nasdag		Enter	symbol, nai	ne or ke's	earch		* 🕞 Elements	s Network Console	» ●8 & 5	i x
UR BUSINESSES QU	OTES MARKE	TS NEWS	INVESTING	ADVANCED	INVESTING P	ERSONAL FINANCE	<html ang="en-us
xmlns:og= nttps://
2008/fbml"></html>	class="wide quotes-re ogp.me/ns#" xmlns:fb="h	ttps://www.facebook.	.com/
FlashQuotes	Select the Tir	meframe:	3 Months				<pre>>#shadow-root (o ><head></head></pre>	pen		- 1
InfoQuotes	_						Q ····▼ <body <="" id="body" p=""></body>	s 30		
STOCK DETAILS	Results for:	3 Month, F	From 15-DEC	-2017 TO 15	-MAR-2018		<pre>subscript>_</pre>	noscript>		
Summary Quote	Date	Open	High	Low	Close / Last	Volume	R if IE 8]</td <td>le Tag Manager (noscrip > <script< td=""><td>ut)></td><td></td></script<></td>	le Tag Manager (noscrip > <script< td=""><td>ut)></td><td></td></script<>	ut)>	
Real-Time Quote	03/15/2018	178.5	180 24	178 0701	178.65	22 676 520	3/ src="https:// [endif?>	www.nasdaq.com/includes	<pre>s/respond.js"></pre>	pt> </td
After Hours Quote	03/14/2018	180.32	180.52	177.91	178.44	29 334 630	a if IE 7]</td <td>> <script www.nasdag.com/includes</script </td> <td>s/zindex.is"><td>t> <1</td></td>	> <script www.nasdag.com/includes</script 	s/zindex.is"> <td>t> <1</td>	t> <1
Pre-market Quote	02/42/2010	100.02	100.02	170.04	170.07	20,004,000	3/ [endif?>		REGIN	
Historical Quote	03/13/2018	102.09	103.5	179.24	119.97	31,404,170	P includes/glob	al-masthead.inc		
Option Chain	03/12/2018	180.29	182.39	180.21	181.72	32,191,070	3/ c1 comscore	tag>		
CHARTS	03/09/2018	177.96	180	177.39	179.98	32,130,360	N html.wide.quotes-resp	onsive.quotes-300-right	ody#body noscript	
Basic Chart	03/08/2018	175.48	177.12	175.07	176.94	24,518,850	M 3/ Styles Event Listener	s DOM Breakpoints Prope	rties	
Interactive Chart	03/07/2018	174.94	175.85	174.27	175.03	31,686,450	Filter	thev .cls + pr	osition 0	
COMPANY NEWS	03/06/2018	177.91	178.25	176.13	176.67	23,273,160	element.style {		margin -	
Company Headlines	03/05/2018	175.21	177.74	174.52	176.82	27,825,140	body {	<style></style>	padding -	
Press Releases	03/02/2018	172.8	176.3	172.45	176.21	38,426,060	background-color prgb(234, 2)	*: 35, 236); 0	710.400 × 2773.500	0
Market Stream	03/01/2018	178 54	179 775	172.66	175	48,706,170	}		-	
STOCK ANALYSIS	02/20/2010	170.06	100 615	170.05	170 12	27 569 090	E Console What	s New		×
Analyst Research	02/28/2018	179.20	180.015	178.00	178.12	37,508,080	top	▼ Filter	Default levels	• •
Guru Analysis	02/27/2018	179.1	180.48	178,16	178.39	38,885,510	1 was preloaded	using link preload but	not used within a f	few
Stock Report	02/26/2018	176.35	179.39	176.21	178.97	37,353,670	wasn't preloade	d for nothing.		
Competitors	02/23/2018	173.67	175.65	173.54	175.5	33,772,050	A The resource ht om/g00/MzoyOjM6	tps://hxyzhas.g00.nasda MTc6MzowOjM%3D/ad?i10c.	<pre>g.c VM135 historica getabdurl=3&ad_chan</pre>	nel=
Stock Consultant	02/22/2018	171.8	173.95	171.71	172.5	30,953,760	1 was preloaded seconds from th	using link preload but window's load event.	not used within a t Please make sure it	few
Stock Comparison	02/21/2018	172.83	174.12	171.01	171.07	37,378,070	wasn't preloade	d for nothing.		
				and the second second						

图 3-4 Chrome 浏览器 Web 开发工具箱

Apple Inc. Comm	ion Stock (A X	+					- 🗆 ×
← → C 企 ☆ 最常访问 ■ × 狐 官方数	 ① ▲ h 5点 ● 新手上路 	nttps://www	/.nasda 址	aq.com/symbol/a	apl/historic	al# 🗱	□ ··· ☆ IN □ 5 年 ● (『 = □ #30版#32
Nasdaq	/	En	ter sym	nbol, name or key	word	/	Search Join NASDAQ or Log In
OUR BUSINESSES	QUOTES V	MARKET		NEWS V INVES		DVANCED INVI	VESTING T PERSONAL FINANCE T MY NASDAQ
STOCK DETAILS	Results for: 3	Month, Fro	om 15-0	DEC-2017 TO 15-	MAR-2018		3/16/2018 11:50:00 AM - Seeking Alpha
Summary Quote	Date	Open	High	Low	Close / Last	Volume	UAL's Bad Week, Apple Magazine, Ford Recall, & Latest Trump Talk 2062018114100 AM, Zerks com
Real-Time Quote	03/15/2018	178.5	180.24	178.0701	178.65	22,676,520	 Qualcomm needs an intermediator more than
After Hours Quote	03/14/2018	180.32	180.52	177.81	178.44	29,334,630	a buyer
Pre-market Quote	03/13/2018	182 59	183.5	179.24	179.97	31 464 170	3/16/2018 11:23:00 AM - Reuters
Historical Quote	0011012010	102.00	100.0	110.24		01,404,170	Plays, Apple's Next A.R.
Option Chain	03/12/2018	180.29	182.39	180.21	181.72	32,191,070	0 3/16/2018 10:51:00 AM - Barron's
CHARTS	03/09/2018	177.96	180	177.39	179.98	32,130,360	 Notable ETF Inflow Detected - QQQ, AAPL,
Basic Chart	03/08/2018	175.48	177.12	175.07	176.94	24,518,850	0 3/16/2018 10:49:23 AM - BNK Invest
Interactive Chart	03/07/2018	174.94	175.85	174.27	175.03	31,686,450	0 More AAPL News & Commentary
□ 空 查看器 □ 控制:	とう 日 調試器	{}样式编	编器	◎性能 ◎内存	三 网络	8存储	日日日の日日
+			0.搜索	R HTML	1	规则	计算后 布局 动圈 字体
htmL					^ = j	盒模型	
<html <br="" class="wide wotes</td><td>-responsive que</td><td>otes-300-rig</td><td>ght">M face</html>	hook com/2008/fhm	1		e			
lang="en-us">	/113# AMILIIS.10	- 110093.774	w.idce	000K.COM/ 2000/ 10M	•	外边距	0
<head> (-) </head>					1.1	边框	E e
	1						増充 0
<pre><loogie lag="" manage<="" pre=""></loogie></pre>	er (noscript)	->				0 0 0	0 1023.2×2000.4 0 0
End Google Tag M</td <td>anager (noscrip</td> <td>ot)></td> <td></td> <td></td> <td></td> <td>1.000</td> <td>0</td>	anager (noscrip	ot)>				1.000	0
</td <td></td> <td></td> <td></td> <td></td> <td>~</td> <td></td> <td>θ</td>					~		θ
<					>		0
html.wide.quotes-responsiv	e.quotes-300-r	. > body#bo	ody		1	1023.2×2800.4	relative

图 3-5 Firefox 浏览器 Web 开发工具箱

3.4.2 示例 2: 解析纳斯达克苹果股票数据

下面通过解析纳斯达克苹果公司股票数据熟悉一下如何使用 BeautifulSoup 库解析 HTML 数据。

在编写代码之前先分析一下纳斯达克股票网页数据,首先需要在浏览器中打开网页,单击 F12 键打开 Web 工具箱,如图 3-6 所示。在打开的 Web 工具箱中单击"查看器"标签,查看 HTML 代码,从中可见所需的数据是包裹在元素的中的,且每行数据都放到一个元素中。



图 3-6 浏览器 Web 工具箱

示例实现代码如下:

// 例实现代码如下:
 import urllib.request
 from bs4 import BeautifulSoup
 # url = 'https://www.nasdaq.com/symbol/aapl/historical#.UWdnJBDMhHk'
 # 本地文件访问
 url = "file:///" + os.path.abspath("./nasdaq-Apple1.html").replace("\\", "/")
 req = urllib.request.Request(url)
 with urllib.request.urlopen(req) as response:

```
data = response.read()
html data = data.decode()
```

<pre>sp = BeautifulSoup(html_data, 'html.pa</pre>	rser')	1
# 返回标签元素		_
<pre>tbody = sp.find('tbody')</pre>		(2)
# 返回标签下所有的元素		
<pre>trlist = tbody.select('tr')</pre>		3
# 保存股票数据列表		
data = []		
for tr in trlist:		4
fields = {}	# 保存一行数据	
# 获得交易日期元素		
<pre>th = tr.find('th')</pre>		
fields['Date'] = th.text	♯ 日期	
# 获得 tr 下的所有 td 元素		
tds = tr.select('td')		
fields['Open'] = tds[0].text	♯ 开盘	
fields['High'] = tds[1].text	# 最高	
fields['Low'] = tds[2].text	# 最低	
fields['Close'] = tds[3].text	# 收盘	
fields['Volume'] = tds[4].text	# 成交量	
data.append(fields)		(5)
		0

print("解析完成。", data)

示例运行后,输出结果如下:

解析完成。[{'Date': '10/04/2022', 'Open': '225.64', 'High': '227.49', 'Low': '223.89', 'Close': '227.01', 'Volume': '34,755,550'}, {'Date': '10/03/2022', 'Open': '218.43', 'High': '220.96', 'Low': '215.132', 'Close': '220.82', 'Volume': '30,352,690'}, {'Date': '10/02/2022', 'Open': '223.06', 'High': '223.58', 'Low': '217.93', 'Close': '218.96', 'Volume': '35,767,260'}, {'Date': '10/01/2022', 'Open': '225.07', 'High': '228.22', 'Low': '224.2', 'Close': '224.59', 'Volume': '36,187,160'}, {'Date': '09/30/2022', '0pen': '220.9', 'High': '224.58', 'Low': '220.79', 'Close': '223.97', 'Volume': '26,318,580'}, {'Date': '09/27/2022', 'Open': '220.54', 'High': '220.96', 'Low': '217.2814', 'Close': '218.82', 'Volume': '25,361,290'}, {'Date': '09/26/2022', 'Open': '220', 'High': '220.94', 'Low': '218.83', 'Close': '219.89', 'Volume': '19,088,310'}, {'Date': '09/25/2022', 'Open': '218.55', 'High': '221.5', 'Low': '217.1402', 'Close': '221.03', 'Volume': '22,481,010'}, {'Date': '09/24/ 2022', 'Open': '221.03', 'High': '222.49', 'Low': '217.19', 'Close': '217.68', 'Volume': '31,434, 370'}, {'Date': '09/23/2022', 'Open': '218.95', 'High': '219.84', 'Low': '217.65', 'Close': '218.72', 'Volume': '19,419,650'}, {'Date': '09/20/2022', 'Open': '221.38', 'High': '222.56', 'Low': '217.473', 'Close': '217.73', 'Volume': '57,977,090'}, {'Date': '09/19/2022', 'Open': '222.01', 'High': '223.76', 'Low': '220.37', 'Close': '220.96', 'Volume': '22,187,880'}, {'Date': '09/18/2022', 'Open': '221.06', 'High': '222.85', 'Low': '219.44', 'Close': '222.77', 'Volume': '25,643,090'}, {'Date' : '09/17/2022', '0pen': '219.96', 'High': '220.82', 'Low': '219.12', 'Close': '220.7', 'Volume': '18,386,470'}, {'Date': '09/16/2022', 'Open': '217.73', 'High': '220.13', 'Low': '217.56', 'Close': '219.9', 'Volume': '21,158,140'}, {'Date': '09/13/2022', 'Open': '220', 'High': '220.79', 'Low': '217.02', 'Close': '218.75', 'Volume': '39,763,300'}, {'Date': '09/12/2022', 'Open': '224.8', 'High': '226.42', 'Low': '222.86', 'Close': '223.085', 'Volume': '32,226,670'}, {'Date': '09/11/2022', 'Open': '218.07', 'High': '223.71', 'Low': '217.73', 'Close': '223.59', 'Volume': '44,289,650'}]

代码解释如下:

42 ┥ 极简ChatGPT: 如何利用AI实现Python自动编程

- •代码第①行使用 BeautifulSoup 构造一个解析器对象 sp,将 HTML 数据作为输入, 并指定解析器为 'html. parser'。
- •代码第②行使用 sp. find('tbody') 查找 HTML 页面中的第一个标签,并将 结果保存在变量 tbody 中。标签通常包含表格数据。
- 代码第③行使用 tbody. select('tr')查找 tbody 标签下的所有
 <hr>
 标签,并将结果保存在列表 trlist 中。每个
 标签表示表格中的一行数据。
 代码第④行遍历 trlist 列表,对于每个
 标签,执行以下操作:
- 使用 tr. find('th')查找当前行中的第一个标签,并将结果保存在变量 th 中。
标签通常用于表示表格中的表头或日期等特殊信息。
- 使用 tr. select('td')查找当前行中的所有标签,并将结果保存在列表tds中。每 个标签表示表格中的一个单元格。
- •代码第⑤行将包含每行数据的字典 fields 添加到列表 data 中。

最后,代码输出解析完成后的股票数据 data。

3.4.3 使用 Selenium 爬取动态网页数据

使用 urllib 爬取数据时经常被服务器反爬技术拦截。服务器有一些办法识别请求是否 来自浏览器。另外,有的数据需要登录系统后才能获得,例如邮箱数据,而且在登录时会有 验证码识别,验证码能够识别出是人工登录系统,还是计算机程序登录系统。试图破解验证 码不是一个好主意,现在的验证码也不是简单的图像,有的会有声音等识别方式。

如果是一个真正的浏览器,那么服务器设置重重"障碍"就不是问题了。Selenium 可以 启动本机浏览器,然后通过程序代码操控它。Selenium 直接操控浏览器,可以返回任何形 式的动态数据。使用 Selenium 操控浏览器的过程中也可以人为干预,例如在登录时,如果 需要输入验证码,则由人工输入,登录成功之后,再由 Selenium 操控浏览器爬取数据。

1. 安装 Selenium

要使用 Selenium 库首先需要安装 Selenium,可通过如下 pip 指令进行安装:

pip install selenium

安装过程如图 3-7 所示。

2. 配置 Selenium

运行 Selenium 需要操作本地浏览器,默认是 Firefox,因此比较推荐安装 Firefox 浏览器,要求 Firefox 浏览器是 55.0 以上版本。由于版本兼容的问题还需要下载浏览器引擎 GeckoDriver,GeckoDriver可以在 https://github.com/mozilla/geckodriver/releases 下载, 根据自己的平台选择对应的版本,不需要安装 GeckoDriver,只需将下载包解压处理就可以了。

然后需要配置环境变量,将 Firefox 浏览器的安装目录和 GeckoDriver 解压目录添加到 系统的 PATH 中,如图 3-8 所示是在 Windows 10 下添加 PATH。

×

III 命令提示符	-		×
Microsoft Windows [版本 10.0.19045.3086]			
(c) Microsoft Corporation。保留所有权利。			
C:\Users\tonysnin_install_selenium			
Collecting selenium			
Downloading selenium-4.10.0-py3-none-any.wh1_(6.7 MB)			
6.7/6.7 MB 4.9 kB/s eta 0:00:00			
Requirement already satisfied: urllib3[socks]<3,>=1.26 in c:\users\tony\appdata\local\programs\python\pythor	311\1	ib∖si	ite
-packages (from selenium) (2.0.3)			
Collecting trio~=0.17 (from selenium)			

Using cached trio-0.22.0-py3-none-any.whl (384 kB)

Collecting trio-websocket~=0.9 (from selenium)

Downloading trio_websocket-0.10.3-py3-none-any.whl (17 kB)

Requirement already satisfied: certifi>=2021.10.8 in c:\users\tony\appdata\local\programs\python\python311\lib\site-pack ages (from selenium) (2023.5.7)

Requirement already satisfied: attrs>=19.2.0 in c:\users\tony\appdata\local\programs\python\python311\lib\site-packages (from trio~=0.17->selenium) (23.1.0)

Collecting sortedcontainers (from trio~=0.17->selenium)

图 3-7 Selenium 安装过程

貫环境变量	
C:\Program Files\Common Files\Oracle\Java\javapath	新建(N)
C:\app\tony\product\21c\dbhomeXE\bin	
C:\Program Files (x86)\VMware\VMware Player\bin\	编辑(E)
%SystemRoot%\system32	
%SystemRoot%	浏览(B)
%SystemRoot%\System32\Wbem	1/355(5)
%SYSTEMROOT%\System32\WindowsPowerShell\v1.0\	開除の
%SYSTEMROOT%\System32\OpenSSH\	(U) ASPERIA
C:\Program Files\Go\bin	
C:\Program Files\Git\cmd	L12/2 (1))
C:\Program Files\dotnet\	上移(0)
C:\Program Files\nodejs\	T 12/0
E:\software\geckodriver-v0.33.0-win32	下移(0)
C:\Program Files\Firefox Developer Edition\	
	编辑文本(T)

图 3-8 添加 PATH

3. Selenium 常用 API

Selenium 操作浏览器主要通过 WebDriver 对象实现, WebDriver 对象提供了操作浏览 器和访问 HTML 代码中数据的函数。

操作浏览器的函数如下:

• refresh(): 刷新网页。

44 🚽 极简ChatGPT: 如何利用AI实现Python自动编程

- back(): 回到上一个页面。
- forward():进入下一个页面。
- close(): 关闭窗口。
- quit():结束浏览器执行。
- get(url): 浏览 URL 所指的网页。

访问 HTML 代码中数据的函数如下:

- find_element(By. ID, id): 通过元素的 id 查找符合条件的第一个元素。
- find_elements(By. ID, id): 通过元素的 id 查找符合条件的所有元素。
- find_element(By.NAME, name): 通过元素名字查找符合条件的第一个元素。
- find_elements(By. NAME, name): 通过元素名字查找符合条件的所有元素。
- find_element(By.LINK_TEXT, link_text):通过链接文本查找符合条件的第一个 元素。
- find_elements(By. LINK_TEXT, link_text): 通过链接文本查找符合条件的所有 元素。
- find_element(By. TAG_NAME, name): 通过标签名查找符合条件的第一个元素。
- find_elements(By. TAG_NAME, name): 通过标签名查找符合条件的所有元素。
- find_element(By. XPATH, xpath):通过 XPath 查找符合条件的第一个元素。
- find_elements(By. XPATH, xpath):通过 XPath 查找符合条件的所有元素。
- find_element(By. CLASS_NAME, name): 通过 CSS 中 class 属性查找符合条件的 第一个元素。
- find_elements(By. CLASS_NAME, name): 通过 CSS 中 class 属性查找符合条件的 所有元素。
- find_element(By. CSS_SELECTOR, css_selector):通过 CSS 中选择器查找符合条 件的第一个元素。
- find_elements(By.CSS_SELECTOR, css_selector): 通过 CSS 中选择器查找符合 条件的所有元素。

3.4.4 示例 3: 爬取搜狐证券贵州茅台股票数据

下面通过爬取搜狐证券贵州茅台股票数据案例,熟悉一下如何使用 Selenium 库爬取和 解析 HTML 数据。

读者如果直接使用 urllib 库是无法直接获取 HTML 数据的,原因是这些数据是同步动态数据。而使用 Selenium 返回这些数据是非常简单的。

在爬取数据之前,还是先分析一下搜狐证券贵州茅台股票的 HTML 数据,主流的浏览 器都提供了 Web 工具箱,找到显示这些数据的 HTML 标签,如图 3-9 所示,在 Web 工具箱 的查看器中,找到显示页面表格对应的 HTML 标签,注意在查看器中选中对应的标签,页 面会将该部分灰色显示。经过查找分析最终找到一个 table 标签,复制它的 id 或 class 属性

值,以备在代码中进行查询。

🔌 贵州茅台:1139.7	9 -0.46% -5.2 🗙	+								-		o x
← → ♂ ŵ	🖸 🔏 q.ste	ock. sohu.com /	′cn/600519/	lshq.shtml			鬷	··· ☆	III\ 🗉	9 4	\$	▶ Ξ
1 Charles Arthorne A	- I - U - U - U - U - U - U - U - U - U	1		w.							2/11-14/	INN AL
最近浏览股 我的自选股 万科A 26.77(-0.41%)	<u>贵州茅台</u> (600519)	1139.79 2020-04-03 15:	*0.46%	-5.21	0	9 实	时行情	● 加入	自选股			
中兴通讯 40.97(-4.05%)	实时行情	成交明细	分	价表	历史行	情	龙	虎榜数据	大宗交易			
	历史行情											
		从 20	019-12-05	table#BIZ hq	history	Searc	ch.tableQ	791.4 × 21	32.8 按周 〇按月	查询		
☆ 贵州茅台	日期	开盘	收盘	涨跌额	涨跌	幅	最低	夏言	成交量(手)	成交金額(万)	换手率
- 27/12/1	累计:	2019-12-05至	2020-04-03	17.46	1.	56%	960.1	1195.5	3383376	3701496	1.21	26.9
行情图表	2020-04-03	1139.03	1139.79	-5.21	-0.	46%	1131.98	1147.96	27618	31484	2.19	0.2
实时行情 推荐	2020-04-02	1104.00	1145.00	42.00	3.	81%	1103.88	1145.00	45206	50951	6.75	0.3
成交明细	2020-04-01	1117.00	1103.00	-8.00	-0.	72%	1103.00	1129.00	33205	37001	8.44	0.2
分价表	2020-03-31	1082.00	1111.00	39.00	3.	64%	1081.80	1115.00	47984	52931	7.06	0.3
< 「 「 「 」 小 香着器 「 」 校制	台口)调试器	↑↓ 网络 {}	样式编辑器		内存	Ð	存储 🕇	无障碍环境	1		a	י גיייה ה
		1¥ P3+4 ()	17.04040400		* *		い 「 」 「 」 」 」	而而	。	改字体	动画	
▼ <div class="</td><td>BIZ innerConter</td><td>t"></div>			^	:ho	v.cls +	▼ 強件念						
♦ <div class<="" p=""></div>	"innerHeaderA"	> ····				元素	• 内联	洗择一	个弹性(Flex)容	器或项目以继续	续.	
▼ <div class<br="">▼<table id<="" td=""><td>= part > ="BIZ_hq_histo</td><td>rySearch" clas:</td><td>s="tableQ"></td><td></td><td></td><td>{ }</td><td></td><td>- 岡格</td><td>· · · · · · · · · · · · · · · · · · ·</td><td></td><td></td><td></td></table></div>	= part > ="BIZ_hq_histo	rySearch" clas:	s="tableQ">			{ }		- 岡格	· · · · · · · · · · · · · · · · · · ·			
<pre> <thead> <thead> </thead></thead></pre>	<pre></pre>					4	.1.css:58	W TATA	F沿右体田 CSS M	秋		
	w v/ coody /				- 1	erC	ontent	<u>щоди</u>	L/反行使用 C33 M	91 <i>H</i>		
<div styl<="" td=""><td>e="padding-top</td><td>:5px">*注:每次</td><td>查询最多显示10</td><td>0条</td></div>	e="padding-top	:5px">*注:每次	查询最多显示10	0条		Cab	width:	▼ 盒楔型				
					~		100%; border:	201		9		
<pre>innerBoard > div.BIZ inner</pre>	erContent > div.p	art ≻ table#BIZ	hg historySeard	ch.tableQ > the	ad >	1	▶ 0;	ma	border	0	_	

图 3-9 Web 工具箱

案例实现代码如下:

from selenium import webdriver from selenium.webdriver.common.by import By

driver = webdriver.Firefox() 2 driver.get('http://q.stock.sohu.com/cn/600519/lshq.shtml') 3 table_element = driver.find_element(By.ID, 'BIZ_hq_historySearch') 4 print(table_element.text) 5 driver.quit() 6

代码解释如下:

- •代码第①行导入 Selenium 库中的 webdriver 模块,该模块提供了用于控制不同浏览 器的驱动程序。
- •代码第②行创建一个 Firefox 浏览器的 WebDriver 实例,将其赋值给变量 driver。这 将启动一个 Firefox 浏览器窗口。
- 代码第③行使用 WebDriver 加载指定的 URL,这里是"http://q.stock.sohu.com/ cn/600519/lshq.shtml",即搜狐股票网站中贵州茅台(股票代码 600519)的历史行情页面。
- •代码第④行使用 WebDriver 的 find_element()函数通过元素的 ID 查找页面上的一

1

个特定元素。这里通过 By. ID 参数指定按照元素的 ID 进行查找, ID 值为'BIZ_hq_ historySearch'。

•代码第⑤行打印找到的元素的文本内容。text 属性返回元素的可见文本。

•代码第⑥行 driver. quit()关闭浏览器窗口并终止 WebDriver 的会话。

3.4.5 示例 4: 使用 Selenium 解析 HTML 数据

Selenium 库不仅可以模拟人工操作 Web 页面,我们也可以利用它的一系列 find_ element()函数进行解析 HTML 数据,使用过程类似于 BeautifulSoup 库,本节来介绍如何 使用 Selenium 库解析搜狐证券贵州茅台股票 HTML 数据。

案例实现代码如下:

```
from selenium import webdriver
from selenium. webdriver. common. by import By
driver = webdriver.Firefox()
driver.get('http://g.stock.sohu.com/cn/600519/lshg.shtml')
table_element = driver.find_element(By.ID, 'BIZ_hq_historySearch') ①
                                                                 (2)
tbody = table_element.find_element(By.TAG_NAME, "tbody")
                                                                 (3)
trlist = tbody.find elements(By.TAG NAME, 'tr')
# 股票数据列表
data = []
for idx, tr in enumerate(trlist):
                                                                 (4)
   if idx == 0:
      # 跳过 table 第一行
      continue
                                                                 (5)
   td list = tr.find elements(By.TAG NAME, "td")
                                                                 (6)
   fields = {}
   fields['Date'] = td_list[0].text
                                         # 日期
   fields['Open'] = td list[1].text
                                       ≠ 开盘
   fields['Close'] = td list[2].text
                                       # 收盘
   fields['Low'] = td list[5].text
                                        # 最低
   fields['High'] = td list[6].text
                                       # 最高
   fields['Volume'] = td list[7].text
                                       # 成交量
   data.append(fields)
```

print(data)
driver.quit()

示例运行后,输出结果如下:

[{'Date': '2023 - 06 - 21', 'Open': '1740.00', 'Close': '1735.83', 'Low': '1735.00', 'High': '1756.60', 'Volume': '17721'}, {'Date': '2023 - 06 - 20', 'Open': '1740.00', 'Close': '1743.46', 'Low': '1735.00', 'High': '1765.00', 'Volume': '20947'}, {'Date': '2023 - 06 - 19', 'Open': '1790.00', 'Close': '1744.00', 'Low': '1738.00', 'High': '1797.95', 'Volume': '31700'}, {'Date': '2023 - 06 - 16', 'Open': '1757.00', 'Close': '1797.69', 'Low': '1750.10', 'High': '1800.00', 'Volume': '37918'}, {'Date': '2023 - 06 -15', 'Open': '1730.34', 'Close': '1715.00', 'Low': '1723.00', 'High': '1755.65', 'Volume': '25223'}, {'Date': '2023 - 06 - 14', 'Open': '1719.00', 'Close': • • •

'1813.74', 'Low': '1783.30', 'High': '1822.01', 'Volume': '23952'}, {'Date': '2023 - 02 - 27', 'Open': '1778.50', 'Close': '1810.41', 'Low': '1775.02', 'High': '1815.00', 'Volume': '22065'}, {'Date': '2023 - 02 - 24', 'Open': '1810.11', 'Close': '1788.00', 'Low': '1782.18', 'High': '1810.19', 'Volume': '24635'}]

代码解释如下:

- •代码第①行使用 WebDriver 在页面中查找具有 ID 为'BIZ_hq_historySearch'的元素,并将其赋值给变量 table_element。这个元素应该是包含历史行情数据的表格。
- •代码第②行在 table_element 元素中查找名为"tbody"的子元素,并将其赋值给变量 tbody。这个操作是为了定位表格中的 tbody 部分,其中包含了行情数据的行。
- 代码第③行在 tbody 元素中查找所有名为"tr"的子元素,返回一个包含这些元素的 列表。这个操作是为了获取每一行行情数据的 tr 元素。
- •代码第④行使用 enumerate()函数遍历 trlist 列表中的每个元素,并为每个元素分配 一个索引 idx 和一个变量 tr,用于迭代行情数据的每一行。
- 代码第⑤行 if idx == 0: continue: 如果索引 idx 等于 0,也就是第一行表头行,就 跳过此次循环,不处理表头行的数据。
- •代码第⑥行在当前行的 tr 元素中查找所有名为"td"的子元素,返回一个包含这些元素的列表。这个操作是为了获取当前行中每个列的数据。

在循环中,每行的各列数据被提取出来,以字典的形式存储在 fields 变量中,并添加到 data 列表中(data 列表包含了每行行情数据的字典)。

最后,代码打印出 data 列表,即历史行情数据,关闭浏览器并终止 WebDriver 的会话。

3.4.6 示例 5: 利用 ChatGPT 辅助解析数据

如果读者编程能力不是很强,也可以利用 ChatGPT 解析 HTML,下面具体介绍。



读者向 ChatGPT 发送如下提问:

根据如下的 HTML 代码,帮我从其中 historical-data__table 表格中解析出数据,注意 不要返回程序代码,你解析返回 CSV 表格数据。

```
<!doctype html >
< html lang = "en">
< head >
< meta charset = "UTF - 8">
< meta name = "Generator" content = "EditPlus®">
< meta name = "Author" content = "">
```

```
< meta name = "Keywords" content = "">
 < meta name = "Description" content = "">
 <title>Document</title>
</head>
< body >
< div id = "quotes content left pnlAJAX">
 < thead class = "historical - data table - headings">
    Date 
      Open 
      High 
      Low 
      Close 
      Volume 
   </thead>
   10/04/2022
     225.64
     227.49
     223.89
     227.01
     34,755,550 
   >09/11/2022
     218.07
     223.71
     217.73
     223.59
     44,289,650 
   </div>
</body>
</html>
执行上述的提问 ChatGPT 回答如下。
```

以下是从给定的 HTML 中解析出的数据,并以 CSV 格式返回的结果:

Date, Open, High, Low, Close, Volume 10/04/2022, 225.64, 227.49, 223.89, 227.01, 34, 755, 550 10/03/2022, 218.43, 220.96, 215.132, 220.82, 30, 352, 690 10/02/2022, 223.06, 223.58, 217.93, 218.96, 35, 767, 260 10/01/2022, 225.07, 228.22, 224.2, 224.59, 36, 187, 160 09/30/2022, 220.9, 224.58, 220.79, 223.97, 26, 318, 580 09/27/2022, 220.54, 220.96, 217.2814, 218.82, 25, 361, 290 09/26/2022, 220, 220.94, 218.83, 219.89, 19, 088, 310 09/25/2022, 218.55, 221.5, 217.1402, 221.03, 22, 481, 010 09/24/2022, 221. 03, 222. 49, 217. 19, 217. 68, 31, 434, 370 09/23/2022, 218. 95, 219. 84, 217. 65, 218. 72, 19, 419, 650 09/20/2022, 221. 38, 222. 56, 217. 473, 217. 73, 57, 977, 090 09/19/2022, 222. 01, 223. 76, 220. 37, 220. 96, 22, 187, 880 09/18/2022, 221. 06, 222. 85, 219. 44, 222. 77, 25, 643, 090 09/17/2022, 219. 96, 220. 82, 219. 12, 220. 7, 18, 386, 470 09/16/2022, 217. 73, 220. 13, 217. 56, 219. 9, 21, 158, 140 09/13/2022, 220, 220. 79, 217. 02, 218. 75, 39, 763, 300 09/12/2022, 224. 8, 226. 42, 222. 86, 223. 085, 32, 226, 670 09/11/2022, 218. 07, 223. 71, 217. 73, 223. 59, 44, 289, 650

这是从表格中提取的日期、开盘价、最高价、最低价、收盘价和成交量的数据。



CSV(Comma Separated Values)是用逗号分隔数据项(也称为字段)的数据交换格式, CSV 主要应用于电子表格和数据库之间的数据交换。CSV 是文本文件,可以使用记事本等 文本编辑工具打开。

3.5 本章总结

本章主要介绍了 ChatGPT 在数据采集中的应用,以及 Python 进行网页爬虫的数据采 集方法。

首先,讨论了 ChatGPT 可以在数据采集的多个方面提供帮助,比如提出采集思路,说明采 集场景,甚至直接提供采集代码。合理利用 ChatGPT 可以极大地提高数据采集的效率。

随后,介绍了数据采集的基本概念,包括其重要性、步骤、技术和工具等。数据采集是后 续数据分析的基础。

接下来,重点介绍了 Python 进行网页爬虫的数据采集。先介绍了网页的静态和动态数据,然后使用 urllib 库爬取了纳斯达克苹果股票数据。为了解析爬取的数据,学习使用了 BeautifulSoup 库的方法。

对于动态网页,介绍了 Selenium 的安装配置和用法,并爬取了搜狐证券的贵州茅台股 票数据。还介绍了如何使用 Selenium 解析 HTML 数据。

最后,以利用 ChatGPT 辅助解析数据为例,展示了 ChatGPT 在数据采集中的强大作用。

通过本章的学习,我们掌握了 Python 进行网页爬虫的数据采集的基本方法,为后续数据分析项目奠定了基础。