

Python Packages Used for Web Scraping:

1. Selenium
2. BeautifulSoup
3. Requests

APIs used: Oxylabs API

It offers real-time scraping functionality for accessing Lazada reviews, retrieving data in a structured format via proxy services. However, it only provides a 7-day free trial, so the authentication credentials must be updated once they expire for the code to continue working.

Link for reference:

1. [GitHub - oxylabs/lazada-scraper: The Lazada Scraper API on Github is a powerful tool designed to retrieve product details, offers, reviews, and seller information from the Lazada e-commerce platform.](#)
2. [Lazada Scraper API | Free Trial \(oxylabs.io\)](#)

URLs for the data source

1. [\[Free Shipping\] KF94 face Mask mask kf94 50pcs malaysia Made in Korea Original 50PCS Washable Cloth Korea k f94 kf95 facemask viral With Design Kf94 Mask Original 50 Pcs Single Facialmask murah【Local Stock】 | Lazada](#)
2. [KF94 Korea Mask 4 LAYER Disposable Earloop Face Mask KF94 Pelitup Muka/ KF94 \(10PCS\) | Lazada](#)
3. [KF94 Korea Earloop Headloop 4 PLY Disposable Earloop Face Mask KF94 Pelitup Muka/ KF94 韩国口罩 | Lazada](#)

Table with brief descriptions of the data fields of the raw data.

Name	The buyer's name who left a review about the product.
Review	The comment or feedback provided by the buyer regarding the product.
SkuInfo	Details about the product category, such as color and quantity (e.g., number of pieces in a mask pack).
Date	The date when the buyer left the review.
StarCount	The number of stars (out of 5) given by the buyer to rate the product.

```
In [ ]: ## Author: Wong Yee En
```

Web Scraping

```
In [2]: # Import the LazadaScrapeWithoutAPI class from the Python file
from lazada_scrape_selenium import LazadaScrapeWithoutAPI

# URL of the Lazada product page to scrape
url = 'https://www.lazada.com.my/products/kf94-korea-mask-4-layer-disposable-earloo

# Initialize the scraper object
scraper = LazadaScrapeWithoutAPI(url, max_pages=3)

# Scrape the reviews
scraper.scrape_reviews()

# Export the reviews to JSON and CSV files
scraper.export_to_json('lazada_reviews.json')
scraper.export_to_csv('lazada_reviews.csv')

# Clear reviews if needed (for a new scraping session)
scraper.clear_reviews()
```

```

1***1 Comfortable fit and stylish design, Easy to wear and fold, Effective four-layer filtration, color_family:White 11 May 2024 5
Mui L. Perfect size for adults, Comfortable fit and stylish design, Convenient earloop design, Variety of colors to choose from, color_family:Black 1 day ago 5
1***5 Waterproof and breathable mask, Variety of colors to choose from, Effective four-layer filtration, Perfect size for adults, color_family:White 1 week ago 5
Mui L. Perfect size for adults, Convenient earloop design, Comfortable fit and stylish design, color_family:White 1 day ago 5
ow K. abit more transparent or thinner material than usual. also smaller. fast delivery. placed on Aug 5, got on Aug 7. color_family:White 07 Aug 2024 4
I***y Variety of colors to choose from,
Perfect size for adults, Easy to wear and fold. color_family:Black 3 weeks ago 5
Chan M. Convenient earloop design, High-quality materials used, color_family:White 1 week ago 5
J***y good value for money
delivery also ok
Thank you Lazada color_family:Black 3 weeks ago 5
NG T. Convenient earloop design, Convenient earloop design, color_family:White 3 weeks ago 5
Coco M. Very cheap at less than 10c each of these mask. The quality is actually quite thin, just normal use, I don't think suitable for patients. By the way it's value buy. color_family:Black 08 May 2024 5
1***1 Secure nose clip for a snug fit, Waterproof and breathable mask, Convenient earloop design, color_family:Black 11 May 2024 5
P***. Easy to wear and fold, Price affordable
Perfect size for adults, color_family:Black 19 Jun 2024 5
zin M. I buy White 8 Not good XXXXX
black colour 🌟🌟🌟🌟🌟🌟🌟🌟
```

```

white Don't buy color_family:White 19 Aug 2023 5
ong T. this face mask good quality and good condition ,
fast delivery. color_family:Black 22 May 2024 5
D***. Goods received in good conditions. Well packed and fast delivery. color_family:White 24 Apr 2024 5
```

API

```
In [3]: # Import the LazadaScrapeWithAPI class
from lazada_scrape_api import LazadaScrapeWithAPI

# List of Lazada product URLs
urls = [
    'https://www.lazada.com.my/products/kf94-korea-earloop-headloop-4-ply-disposable',
    'https://www.lazada.com.my/products/kf94-korea-mask-4-layer-disposable-earloop'
]

# Authentication credentials for OxyLabs API
auth_credentials = ('YeeEn_J1wov', 'Yeenlenglui00~')

# Initialize the scraper object
scraper = LazadaScrapeWithAPI(auth_credentials)

# Scrape reviews from the list of URLs
```

```
scraper.scrape_reviews(urls)

# Export the scraped data to JSON and CSV
scraper.export_to_json('reviews_api_example.json')
scraper.export_to_csv('reviews_api_example.csv')

# Clear reviews
scraper.clear_reviews()
```

In []:

```

1 ## Author: Wong Yee En
2 import requests
3 from bs4 import BeautifulSoup
4 import json
5 import csv
6
7 class LazadaScrapeWithAPI:
8     def __init__(self, auth_credentials):
9         """
10             Initialize the LazadaScrapeWithAPI class with authentication credentials for the
11             Oxylabs API.
12         """
13         self.auth_credentials = auth_credentials
14         self.all_reviews = []
15
16     def scrape_reviews(self, urls):
17         """
18             Scrape reviews from a list of Lazada product URLs using the Oxylabs API.
19
20             Args:
21                 urls (list): List of Lazada product URLs to scrape.
22
23             Returns:
24                 None
25         """
26         for url in urls:
27             # Set up proxy with requests
28             payload = {
29                 'source': 'universal',
30                 'url': url,
31                 'render': 'html',
32                 'user_agent_type': 'desktop',
33                 'context': [{'key': 'follow_redirects', 'value': True}],
34             }
35
36             response = requests.post(
37                 'https://realtime.oxylabs.io/v1/queries',
38                 auth=self.auth_credentials,
39                 json=payload,
40             )
41
42             html_content = response.json()['results'][0]['content']
43
44             # Parse the HTML content with BeautifulSoup
45             soup = BeautifulSoup(html_content, 'html.parser')
46
47             # Extract reviews for each URL
48             nameContainers = soup.findAll('div', class_='middle')
49             containers = soup.findAll('div', class_='item-content')
50             dateContainers = soup.findAll('div', class_='top')
51
52             for nameContainer, container, dateContainer in zip(nameContainers, containers, dateContainers):
53                 name = nameContainer.find('span')
54                 review = container.find('div', class_='content')
55                 skuInfo = container.find('div', class_='skuInfo')

```

9/7/24, 3:48 PM

lazada_scrape_api.py - Jupyter Text Editor

```
55         date = dateContainer.find('span', class_='title right')
56         stars = dateContainer.find('div', class_='container-star starCtn'
57 left').find_all('img', class_='star',
58 src='//img.lazcdn.com/g/tps/tfs/TB19ZvEgfDH8KJjy1XcXXcpdXXa-64-64.png')
59         star_count = len(stars)
60
61     if name and review and skuInfo and date and stars:
62         name_text = name.get_text(strip=True)
63         review_text = review.get_text(strip=True)
64         skuInfo_text = skuInfo.get_text(strip=True)
65         date_text = date.get_text(strip=True)
66         self.all_reviews.append((name_text, review_text, skuInfo_text,
67 date_text, star_count))
68
69     def export_to_json(self, file_name):
70         """
71             Export the scraped data to a JSON file.
72
73         Args:
74             file_name (str): The name of the JSON file to export the data.
75
76         Returns:
77             None
78         """
79         with open(file_name, 'w', encoding='utf-8') as f:
80             json.dump(self.all_reviews, f, ensure_ascii=False, indent=4)
81
82     def export_to_csv(self, file_name):
83         """
84             Export the scraped data to a CSV file.
85
86         Args:
87             file_name (str): The name of the CSV file to export the data.
88
89         Returns:
90             None
91         """
92         with open(file_name, 'w', newline='', encoding='utf-8') as f:
93             writer = csv.writer(f)
94             writer.writerow(['Name', 'Review', 'SkuInfo', 'Date', 'StarCount'])
95             for name, review, skuInfo, date, star_count in self.all_reviews:
96                 writer.writerow([name, review, skuInfo, date, star_count])
97
98     def clear_reviews(self):
99         """
100            Clear the stored reviews.
101
102        Returns:
103            None
104        """
105        self.all_reviews = []
```

```

1 ## Author: Wong Yee En
2 import time
3 from selenium import webdriver
4 from bs4 import BeautifulSoup
5 import json
6 import csv
7 from selenium.webdriver.common.by import By
8 from selenium.webdriver.common.keys import Keys
9 from selenium.webdriver.support.ui import WebDriverWait
10 from selenium.webdriver.support import expected_conditions as EC
11
12 class LazadaScrapeWithoutAPI:
13     def __init__(self, url, max_pages=5):
14         """
15             Initialize the LazadaScrapeWithoutAPI class.
16
17         Args:
18             url (str): The URL of the Lazada product page to scrape.
19             max_pages (int): The maximum number of pages to scrape reviews from.
20         """
21         self.url = url
22         self.max_pages = max_pages
23         self.reviews = []
24         self.driver = None
25
26     def initialize_driver(self):
27         """
28             Initialize the Selenium WebDriver with Chrome options.
29         """
30         options = webdriver.ChromeOptions()
31         options.add_argument("--start-maximized")
32         self.driver = webdriver.Chrome(options=options)
33         self.driver.get(self.url)
34         body = self.driver.find_element(By.TAG_NAME, 'body')
35         for _ in range(10): # Adjust the range for more/less increments
36             body.send_keys(Keys.ARROW_DOWN) # Scroll a small amount down
37             body.send_keys(Keys.ARROW_DOWN)
38             body.send_keys(Keys.ARROW_DOWN)
39             time.sleep(2) # Adjust sleep duration for slower or faster scrolling
40
41     def scrape_reviews(self):
42         """
43             Scrape the reviews from the Lazada product page.
44         """
45         self.initialize_driver()
46
47         for i in range(0, self.max_pages):
48             soup = BeautifulSoup(self.driver.page_source, 'html.parser')
49
50             nameContainers = soup.findAll('div', class_='middle')
51             containers = soup.findAll('div', class_='item-content')
52             dateContainers = soup.findAll('div', class_='top')
53
54             for nameContainer, container, dateContainer in zip(nameContainers, containers,
55             dateContainers):
55                 name = nameContainer.find('span')

```

9/7/24, 3:48 PM

```
lazada_scrape_selenium.py - Jupyter Text Editor

56     review = container.find('div', class_='content')
57     skuInfo = container.find('div', class_='skuInfo')
58     date = dateContainer.find('span', class_='title right')
59     stars = dateContainer.find('div', class_='container-star starCtn'
59     left').find_all('img', class_='star',
60     src='//img.lazcdn.com/g/tps/tfs/TB19ZvEgfDH8KJjy1XcXXcpdXXa-64-64.png')
60     star_count = len(stars)
61
62     if name and review and skuInfo and date and stars:
63         name_text = name.get_text(strip=True)
64         review_text = review.get_text(strip=True)
65         skuInfo_text = skuInfo.get_text(strip=True)
66         date_text = date.get_text(strip=True)
67         self.reviews.append((name_text, review_text, skuInfo_text, date_text,
67     star_count))
68         print(name_text, review_text, skuInfo_text, date_text, star_count)
69
70     time.sleep(2)
71
72     # Try to click the next button to go to the next page
73     try:
74         next_button = self.driver.find_element(By.CSS_SELECTOR, "button[class='next
74     btn next-btn-normal next-btn-medium next-pagination-item next']")
75         next_button.click()
76     except:
77         print("No more pages.")
78         break
79
80     time.sleep(3)
81
82     self.driver.quit()
83
84     def export_to_json(self, file_name):
85         """
86             Export the scraped reviews to a JSON file.
87
88             Args:
89                 file_name (str): The name of the JSON file to export the data.
89             """
90
91         with open(file_name, 'w', encoding='utf-8') as f:
92             json.dump(self.reviews, f, ensure_ascii=False, indent=4)
93
94     def export_to_csv(self, file_name):
95         """
96             Export the scraped reviews to a CSV file.
97
98             Args:
99                 file_name (str): The name of the CSV file to export the data.
99             """
100
101        with open(file_name, 'w', newline='', encoding='utf-8') as f:
102            writer = csv.writer(f)
103            writer.writerow(['Name', 'Review', 'SkuInfo', 'Date', 'StarCount'])
104            for name, review, skuInfo, date, star_count in self.reviews:
105                writer.writerow([name, review, skuInfo, date, star_count])
106
107    def clear_reviews(self):
108        """
```

9/7/24, 3:48 PM

lazada_scrape_selenium.py - Jupyter Text Editor

```
109     Clear the stored reviews.  
110     """  
111     self.reviews = []  
112
```

Data Ingestion

```
In [14]: ## Author: Wong Yee En
```

```
In [4]: from pyspark.sql import SparkSession
spark = SparkSession\
    .builder\
    .appName("Ingestion")\
    .getOrCreate()
```

24/09/07 00:53:36 WARN SparkSession: Using an existing Spark session; only runtime SQL configurations will take effect.

```
In [5]: sc = spark.sparkContext
sc.addFile("../de_classes/data_storage/hadoop_file_handler.py")
```

24/09/07 00:53:36 WARN SparkContext: The path ../de_classes/data_storage/hadoop_file_handler.py has been added already. Overwriting of added paths is not supported in the current version.

```
In [6]: # Import the HadoopFileHandler class
from hadoop_file_handler import HadoopFileHandler

# Create an instance of HadoopFileHandler
handler = HadoopFileHandler()

df = handler.read_csv("data/raw/" + "*.csv")

output_path = "data/merged/merged_reviews.csv"
handler.write_csv(df, output_path)
```

24/09/07 00:53:37 WARN SparkSession: Using an existing Spark session; only runtime SQL configurations will take effect.

```
In [7]: df = handler.read_csv('data/merged/merged_reviews.csv')
```

```
In [8]: total_rows = df.count()
print(f"Total number of rows: {total_rows}")
```

Total number of rows: 5767

```
In [9]: df.show(truncate=True)
```

Name	Review	SkuInfo	Date	StarCount
1***1 Comfortable fit a...	color_family:White	11 May 2024	5	
Coco M. Very cheap at les...	color_family:Black	08 May 2024	5	
1***1 Secure nose clip ...	color_family:Black	11 May 2024	5	
P***. Easy to wear and ...	color_family:Black	19 Jun 2024	5	
ong T. this face mask go...	color_family:Black	22 May 2024	5	
zin M. I buy White 8 Not...	color_family:White	19 Aug 2023	5	
D***. Goods received in...	color_family:White	24 Apr 2024	5	
G***n overall ok. U get...	color_family:Black	17 Apr 2024	5	
nur N. CANTIKKK GILE 🌸 ...	Color Family:Pink	15 Feb 2022	5	
Lim C. Good packing, rec...	color_family:Navy...	15 Dec 2023	5	
Joanne L. Goods received in...	color_family:Black	08 Dec 2023	4	
Vivi High-quality mate...	color_family:White	1 day ago	5	
susilawati maaf boss baru se...	color_family:Black	20 Oct 2023	5	
digitalvault So far good, I wi...	color_family:Black	2 weeks ago	5	
****2 Pretty decent mas...	color_family:Black	2 weeks ago	5	
R***. High-quality mate...	color_family:White	1 week ago	5	
Suki S. Received in short...	color_family:Red	02 Oct 2023	5	
Suki S. Received in short...	color_family:Purple	02 Oct 2023	5	
Suki S. Received in short...	color_family:Pink	02 Oct 2023	5	
Jay L. High-quality mate...	color_family:White	29 Jun 2024	5	

only showing top 20 rows

In [10]: `spark.stop()`