

Comparative market analysis of Apple products in US and Switzerland

Final report

Group 51

Xuekai Li

Dmytro Rudyka

Sasa Ljubisavljevic

HSLU

Master in Applied Information and Data Science

Module: Data collection, Integration and Pre-processing

April 2023, Lucerne

Table of contents

Table of contents.....	2
Introduction.....	3
Subject area.....	3
Why this topic.....	3
Sources.....	4
Research questions.....	4
Workflow diagram.....	5
Description of students' tasks.....	5
Student A: Xuekai Li.....	5
Student B: Dmytro Rudyka.....	7
Student C: Sasa Ljubisavljevic.....	8
Description of the process.....	8
Student A: Xuekai Li — www.fust.ch.....	8
Extract.....	8
Transform.....	9
Uploading.....	10
Student B: Dmytro Rudyka — www.abt.com.....	10
Extraction.....	10
Impurifying.....	12
Cleaning, Transformation and Uploading.....	12
Student C: Sasa Ljubisavljevic — www.cyberport.de.....	12
Extract.....	12
Impurifying and cleaning.....	13
Transform.....	13
Uploading.....	13
Merging.....	13
Analysis and answers to questions.....	14
Limitations.....	17
Tools used.....	17
Lesson learned / self-reflection.....	18

Introduction

Subject area

Understanding price differences for goods and services in different markets is crucial for businesses and customers. This is particularly relevant in the technology industry, where Apple products are highly popular and often carry a premium price tag. In this project, we aimed to compare the prices of Apple products between the US, Switzerland and Germany, important markets for Apple. To achieve this we used data integration, preparation and processing techniques to scrape and extract data from two typical retailers in each market, ABT.com in the US, FUST.ch in Switzerland and cyberport.de in Germany.

We have transformed and loaded this data into a suitable data storage system and performed exploratory data analysis to answer to the stated questions. We gained insights into the pricing of Apple products in different markets and are hoping to provide valuable information.

Why this topic

Everyone on our team loves Apple. Well, not everyone. But no one can deny the huge influence of this brand on the entire industry. Statistically, every second phone in Switzerland is an iPhone, and two of the team members have a MacBook. And we have all often heard that the cheapest place to buy it is in the USA. But how true is that? How much difference is there for different types of gadgets? This is very interesting to find out! That's why we decided to do this.

It's also interesting how much of the difference is influenced by the popularity of the products, which assessed using ratings, quantity and quality. We understand that when looking at products of the same brand a big contribution to the price is made by the manufacturer. But we investigated at independent shops to reduce this effect.

Sources

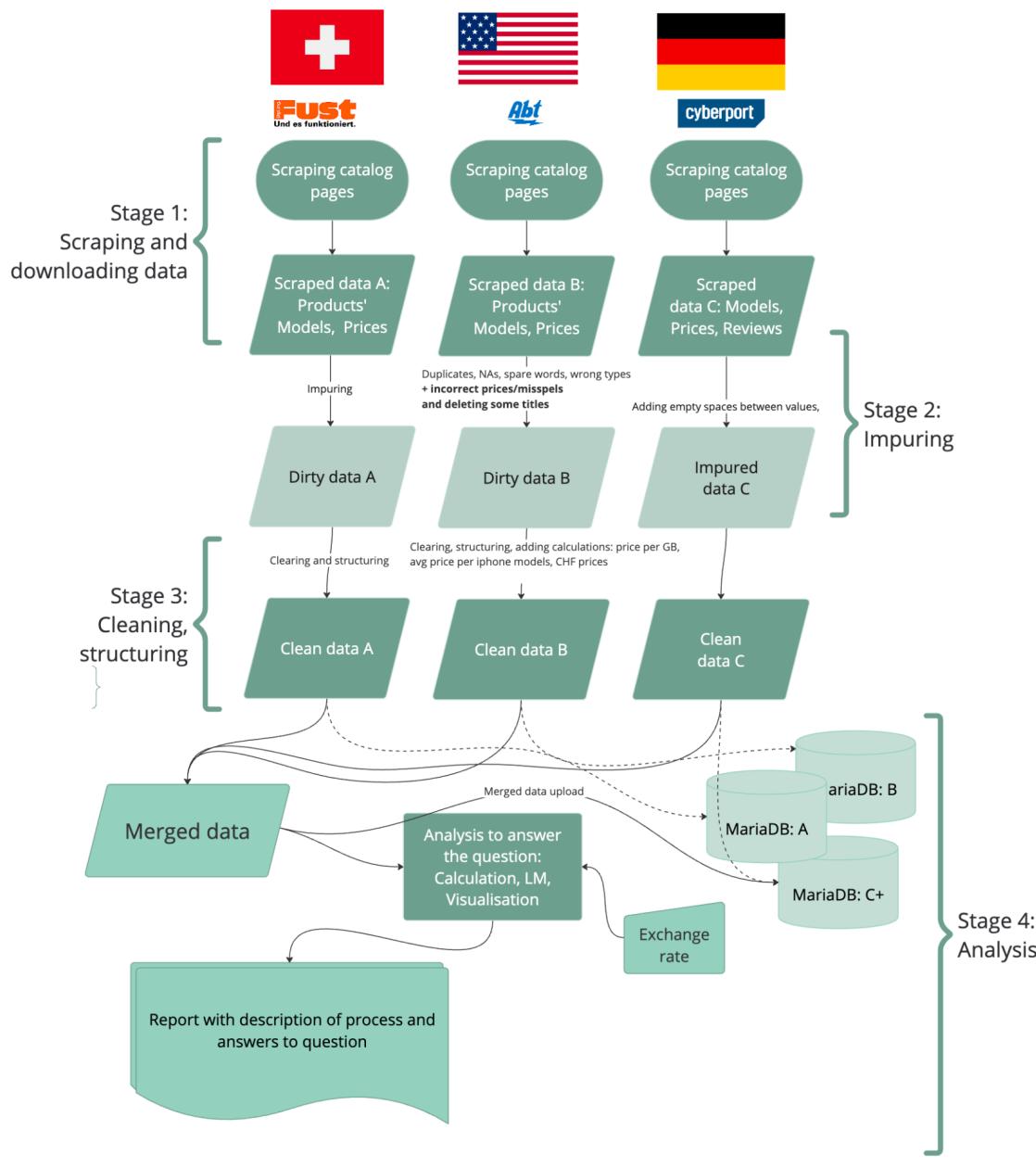
- Source A: scraping from www.fust.ch
- Source B: scraping from www.abt.com
- Source C: scraping from www.cyberport.de

Research questions

By analyzing the selected markets, Switzerland, US, and Germany we aim to answer the following questions:

1. Which country has higher prices for Apple products, the USA, Switzerland or Germany? And by how much?
2. Between Apple mobile phones, laptops and tablets, which category has the biggest price difference?
3. Is there a correlation between the popularity (in terms of number of reviews and overall rating) and the price difference of Apple products in the country with the lowest and highest prices?

Workflow diagram



Description of students' tasks

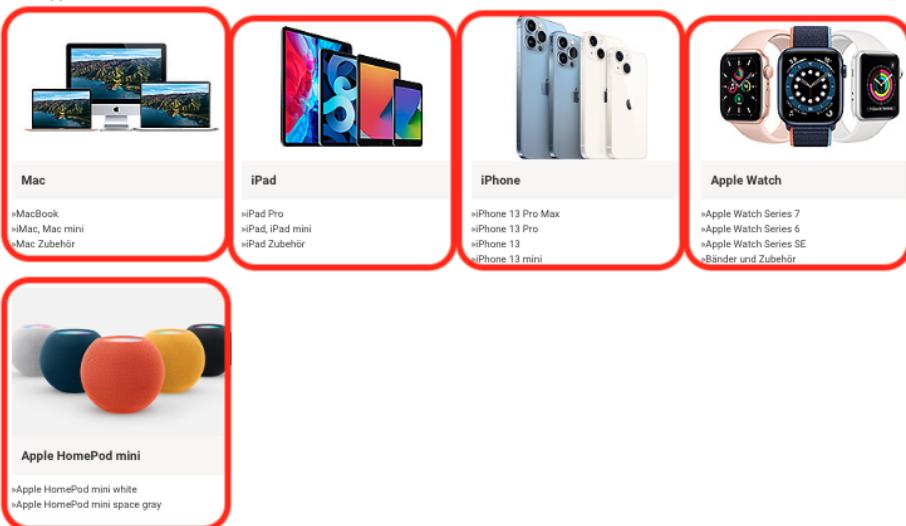
Student A: Xuekai Li

Xuekai scraped data from web source A: <https://www.fust.ch/de/marken/apple.html>

1. category list:

Apple - Markenseite

Das Apple Sortiment



2. item details.

The screenshot shows the Fust website interface. It includes:

- Fust** logo: Und es funktioniert.
- Search bar: suchen...
- Navigation menu: Filialen, Merkliste, Vergleichsliste, Warenkorb, Login
- Category menu: Haushalt, Küchengeräte, PC · Tablet · Handy, TV · Foto · Gaming, Küchen · Badezimmer, Aktionen, Service

Home / PC · Tablet · Handy / Smartphone / Apple iPhone / iPhone 14 Pro

iPhone 14 Pro

The screenshot shows the iPhone 14 Pro product listing on the Fust website. It includes:

Marke	Preis	Arbeitsspeicher in GB	Intern Speicher
Auswählen	von 1059 bis 1699 CHF	von 0 bis 0	von 1 bis 512 GB

Alle Filter anzeigen

Sortierung: beliebteste

1 - 16 von 16

Four iPhone 14 Pro variants are highlighted with red boxes:

- 1'059.-** iPhone 14 Pro, Apple iPhone 14 Pro, 128 GB, Deep Purple, 6,1", 48 MP, 5G (0.0)
- 1'059.-** iPhone 14 Pro, Apple iPhone 14 Pro, 128 GB, Silver, 6,1", 48 MP, 5G (0.0)
- 1'059.-** iPhone 14 Pro, Apple iPhone 14 Pro, 128 GB, Space Black, 6,1", 48 MP, 5G (0.0)
- 1'159.-** iPhone 14 Pro, Apple iPhone 14 Pro, 256 GB, Space Black, 6,1", 48 MP, 5G (0.0)

Each variant has a checkbox for Schneller Gratis-Versand, a delivery status (sofort lieferbar), and availability in 18 or 13 filialen. There are also Merken and Vergleich buttons.

- Eligibility check: www.fust.ch/robots.txt

```
# robots.txt for https://www.fust.ch/
```

```
User-agent:*
Disallow:/admin/
Disallow:/*jsessionid
```

- target attribute: category, model, price, ratings
- made dirty/cleans/enriches the data with Python/Pandas.
- uploaded last stage of personal data before merge to personal Maria DB.

Student B: Dmytro Rudyka

Dmytro scraped catalog pages of Apple products of US online shop abt.com, Aple products are here: <https://www.abt.com/brand/apple>



Eligibility check: After investigating the website [robots.txt](#) it was not found any prohibitions for scraping needed pages:

```
User-agent: *
Disallow: /includes/
Disallow: /mobile/order/*
Disallow: /mobile/cart.php*
Disallow: /mobile/email_friend.php*
Disallow: /mobile/search.php?*
Disallow: /mobile/search_proxy.php*
Disallow: /resources/pages/maintenance.html
Disallow: /resources/pages/category_print.php*
Disallow: /resources/pages/search_proxy.php*
Disallow: /resources/pages/search.php*
Disallow: /list/*
Disallow: /category/*
Disallow: /page/*
Disallow: /resources/pages/capture_zipcode.php*
Disallow: /resources/pages/cart.php*
```

It was not scraped any of disallowed URLs or folders above.

Scraping procedure:

1. Brand catalog with categories is starting point which has the links to the listing with list of the products with the categories

The screenshot shows the Abt.com homepage. At the top, there's a banner for the "87TH ANNIVERSARY SALE" with the text "Save throughout the site" and a "SHOP NOW" button. Below the banner, there's a navigation bar with links for "Brands", "Blog", "Gifts", "Learn", "Videos", "Track Order", "Installation & Services", "SEARCH", "Deliver to My Location", "Account", and "Hello, Sign In". There's also a link to "Call us at 800-860-3577, chat or email". The main content area features a grid of seven Apple product categories: Mac, iPad, iPhone, Apple Watch, Apple TV, AirPods & Earbuds, and AirTag. Each category has a small image and the product name above it. The "iPhone" category is highlighted with a red box.

After that (for some categories through interim page with subcategories) was got the list of pages with the products, for example:

The screenshot shows a product listing page for Apple iPhone 14 models. On the left, there are filters for "Availability" (In Stock), "Brand" (Apple), "Color" (Black, Purple, Red, Blue, Starlight), and "Price" (Less than \$799, \$799 to \$899, \$899 to \$1099, \$1099 to \$1199, \$1199 and up). The main area displays four iPhone 14 models: Gold, Midnight, Blue, and Red. Each model has a small image, the price (\$1,199.99, \$999.99, \$999.99, \$799.99), and a "CALL TO ORDER" button. The "Gold" model is highlighted with a red box. The page is numbered 1 to 5, with "Previous" and "Next" buttons.

2. Went through pages and scrape following **List of attributes**:

- Scraped URL

- Category
- Title with model ID
- Price
- URL of product
- Number of reviews
- Rating

3. Impure

4. Cleaned, extracted and uploaded to MariaDB

Student C: Sasa Ljubisavljevic

Initial plan was to download the Amazon data from BrightData.com. But that could not be achieved in four weeks. After long and slow communication with the vendors and consultants there, we were never able to get an adequate response.

So we decided to scrape one more website with the same data but in new country, Germany. It was www.cyberport.de. The whole procedure is described below.

Description of the process

The processes of scraping each site and processing the information to answer the questions are described below.

Student A: Xuekai Li — www.fust.ch

Extract

After I filtered out the product by apple brand and investigated the website(screenshot below), I found out the page is dynamic, so choose selenium to scrape the data.

1. Get the total items number and calculate how many clicks needed to load all products.

2. Because the web site is dynamic page, choose Selenium to click the button.

After loading all products, I can use selenium to get the product content of each product, and use a for loop to extract the key attributes(screenshot below).

div.product__content.d-11
ex.flex-column.flex-grow 325.97 x 355.98
-1.order-2

```

    <div class="product w-100 p-4 bg-white rounded-slop border-bottom border-shade-box d-flex flex-column flex-grow-1 position-relative">
      <flex>
        <figure class="my-5 mx-auto order-1 position-static product_overview-img aspect-ratio-1/1">...</figure>
        <div class="product__content d-flex flex-column flex-grow-1 order-2">...</div>
      </flex>
    </div>
    </li>
  <li class="product__list d-flex justify-content-between">
    ...
  </li>

```

Finally I transpose the output as panda's data frame, and save as 'stage1.csv' file. From the screenshot below we know it is dirty and price and reviews need to normalize. you will see more details in the next phase.

brand		title	reviews	rating	price
806	Apple	iPad Pro Wi-Fi + 5G 2021 [12.9", 1 TB, Space G...	(0)	0.0	1'999.-
807	Apple	iPad Pro Wi-Fi 2022 [11", 2 TB, Silver, MNXN3T...	(0)	0.0	2'139.-
808	Apple	iPad Pro Wi-Fi 2022 [11", 2 TB, Space Gray, MN...	(0)	0.0	2'139.-
809	Apple	iPad Pro Wi-Fi + 5G 2022 [11", 2 TB, Silver, M...	(0)	0.0	2'309.-
810	Apple	iPad Pro Wi-Fi + 5G 2022 [11", 2 TB, Space Gra...	(0)	0.0	2'309.-
811	Apple	iPad Pro Wi-Fi 2022 [12.9", 2 TB, Silver, MNY0...	(0)	0.0	2'458.-
812	Apple	iPad Pro Wi-Fi 2022 [12.9", 2 TB, Space Gray, ...	(0)	0.0	2'458.-
813	Apple	iPad Pro Wi-Fi + 5G 2022 [12.9", 2 TB, Silver,...	(0)	0.0	2'624.-
814	Apple	Mac Studio 2022 [M1 Ultra Chip, 64 GB RAM, 1 T...	(0)	0.0	4'369.-
815	Apple	Gift Card CHF 30.00	(0)	0.0	33.90

Transform

On 'fust' website didn't provide the product number for 'iPhone', so I need to enrich additional information from the title. That information will be used in the merge phase. I will follow the steps:

Enrichment

1. extract the category from the title.
2. extract the product_number, model, storage,(chip/size/color) from each category.

Cleaning:

1. drop the other category(accessory).
2. drop the outliers

Transformation:

1. normalize the reviews.
2. normalize the price, from string to numeric, and remove "" and "-" from "2'199.-".
3. transform the storage format(from "64 GB" to "64GB").

the cleaned data will look like this:

	brand	title	reviews	rating	price	category	storage	model	sub	prod_n
798	Apple	iPhone 13 Pro Max - 1 TB, Silver, 6.7", 12 MP, 5G	0	0.0	1599.0	iPhone	1TB	13 Pro Max	Silver	-1
799	Apple	iPad Pro Wi-Fi 2022 [11", 1 TB, Silver, MNXL3TY/A]	0	0.0	1689.0	iPad	1TB	Pro	11"	MNXL3TY/A
800	Apple	iPhone 13 Pro - 1 TB, Alpine Green, 6.1", 12 M...	0	0.0	1737.0	iPhone	1TB	13 Pro	Alpine Green	-1
801	Apple	iPad Pro Wi-Fi + 5G 2022 [12.9", 512 GB, Silve...	0	0.0	1744.0	iPad	512GB	Pro	12.9"	MP233TY/A
802	Apple	iMac 2020 [27", Intel Core i5, 16 GB RAM, 256 ...	0	0.0	1749.0	iMac	256GB	2020	i5	-1
803	Apple	iPhone 13 Pro Max - 1 TB, Gold, 6.7", 12 MP, 5G	0	0.0	1749.0	iPhone	1TB	13 Pro Max	Gold	-1
804	Apple	iPhone 13 Pro Max - 1 TB, Alpine Green, 6.7", ...	0	0.0	1837.0	iPhone	1TB	13 Pro Max	Alpine Green	-1
805	Apple	iPad Pro Wi-Fi + 5G 2022 [11", 1 TB, Silver, M...	0	0.0	1869.0	iPad	1TB	Pro	11"	MNYK3TY/A
806	Apple	iPad Pro Wi-Fi + 5G 2021 [12.9", 1 TB, Space G...	0	0.0	1999.0	iPad	1TB	Pro	12.9"	MHRA3TY/A
807	Apple	iPad Pro Wi-Fi 2022 [11", 2 TB, Silver, MNXN3TY...	0	0.0	2139.0	iPad	2TB	Pro	11"	MNXN3TY/A

Uploading

In the final phase, I will use python to upload the data from dataframe. after I connected the 'CIP' database in mariADB, I can create the table as 'fust_stage3'.

```
MariaDB [sys]> use CIP;
Reading table information for
You can turn off this feature
Database changed
MariaDB [CIP]> show tables;
+-----+
| Tables_in_CIP |
+-----+
| fust_stage3   |
| tbl_kontakt  |
+-----+
```

MariaDB [CIP]> show columns from fust_stage3;					
Field	Type	Null	Key	Default	Extra
id	int(11)	NO	PRI	NULL	auto_increment
brand	varchar(255)	YES		NULL	
title	varchar(255)	YES		NULL	
reviews	int(11)	YES		NULL	
rating	float	YES		NULL	
price	float	YES		NULL	
category	varchar(255)	YES		NULL	
storage	varchar(255)	YES		NULL	
model	varchar(255)	YES		NULL	
sub	varchar(255)	YES		NULL	
prod_n	varchar(255)	YES		NULL	

Finally use a for loop to insert each row from the dataframe. and the result is:

	id	brand	title	reviews	rating	price	category	storage	model	sub	prod_n
1	Apple	Watch SE, 44mm, Midnight	0	0	294.9	Watch	44mm	SE	Midnight	-1	
2	Apple	iPad Wi-Fi 2021 [10.2", 64 GB, Silver, MK2L3TY/A]	0	0	298	iPad	64GB	2021	10.2"	MK2L3TY/A	
3	Apple	iPad Wi-Fi 2021 [10.2", 64 GB, Space Gray, MK2K3TY/A]	0	0	298	iPad	64GB	2021	10.2"	MK2K3TY/A	
4	Apple	Watch Series 8, 41mm, Midnight	0	0	423	Watch	41mm	Series 8	Midnight	1	
5	Apple	Watch Series 8, 45mm, Silver	0	0	449	Watch	45mm	Series 8	Silver	-1	
6	Apple	iPad Wi-Fi 2022 [10.9", 64 GB, Silver, MPQ03TY/A]	0	0	466	iPad	64GB	2022	10.9"	MPQ03TY/A	
7	Apple	iPad Wi-Fi + 4G 2021 [10.2", 64 GB, Space Gray, MK473TY/A]	0	0	479	iPad	64GB	2021	10.2"	MK473TY/A	
8	Apple	iPad Wi-Fi 2021 [10.2", 256 GB, Space Gray, MK2N3TY/A]	0	0	498	iPad	256GB	2021	10.2"	MK2N3TY/A	
9	Apple	Watch Series 8, 41mm, Cellular, Midnight	0	0	519	Watch	41mm	Series 8	Midnight	-1	
10	Apple	iPad mini Wi-Fi 2021 [8.3", 64 GB, Space Gray, MK7M3TY/A]	0	0	526	iPad	64GB	mini	-1	MK7M3TY/A	
11	Apple	iPad Air Wi-Fi 2022 [10.9", 64 GB, Blue, MM9E3TY/A]	0	0	534	iPad	64GB	Air	10.9"	MM9E3TY/A	
12	Apple	iPad A1r Wi-Fi 2022 [10.9", 64 GB, Space Gray, MM9C3TY/A]	0	0	534	iPad	64GB	Air	10.9"	MM9C3TY/A	
13	Apple	iPad A1r Wi-Fi 2022 [10.9", 64 GB, Purple, MME23TY/A]	0	0	534	iPad	64GB	Air	10.9"	MME23TY/A	
14	Apple	iPad Wi-Fi 2022 [10.9", 256 GB, Silver, MPQ83TY/A]	0	0	648	iPad	256GB	2022	10.9"	MPQ83TY/A	
15	Apple	iPad A1r Wi-Fi 2022 [10.9", 256 GB, Purple, MME63TY/A]	0	0	773	iPad	256GB	Air	10.9"	MME63TY/A	

Student B: Dmytro Rudyka — www.abt.com

Extraction

Started from the main page of Apple <https://www.abt.com/brand/apple>. From which URLs to scrape was gathered.

Dimensions: Responsive ▾ 522 x 667 100% Custom

The screenshot shows the 'Browse Apple Categories' section of the Abt website. On the left, there's a sidebar with icons for 'Computers & Tablets', 'Headphones', and 'Phones'. The 'Computers & Tablets' item is highlighted with a red box. On the right, the DOM tree for the 'category-widget' is displayed, specifically focusing on the 'cat_list_link' and 'subcat_header' elements for the 'Computers & Tablets' category.

After gathering URLs from the main page and pages with subcategories was founded proper classes to scrape products properties:

The screenshot shows a product page for 'Apple iPhones'. It features a search bar, a 'MEMORIAL DAY SALE' banner, and a 'Helpful Resources' dropdown. The main content displays a list of products, with one item for an 'Apple 128GB Deep Purple iPhone 14 Pro Cellular Phone - MQ0E3LL/A & 6514D' highlighted with a red box. The DOM tree on the right shows the detailed structure of this product listing, including its title, price (\$999.99), and product ID (TPHONF14PROPIRPI F-17RGR).

Scrape the products within all category pages the following function with recursion was applied:

```

1 # function to scrape products and next url altogether, and write the result at once in CSV
2
3 def scrape_prods_to_csv(url, category): #category will be stored in CSV as well, because it will be used on next s
4     print(category)
5     with open('../data/abt_rudyka_stagel.csv', 'a', newline='') as csvfile: #method is in order to just mini
6         writer = csv.writer(csvfile, delimiter=';', quotechar='', quoting=csv.QUOTE_MINIMAL) #quoting= just in ca
7         sleep = random.randint(1,5)
8         print('Waiting for', sleep, 'and then scrape', url) # to be informed of the process
9         time.sleep(sleep)
10        r = requests.get(url)
11        doc = BeautifulSoup(r.text, "html.parser")
12        products = doc.select(".category_item_container")
13        for product in products:
14            title = product.select_one('.cl_title').text.strip() #title of the product
15            price = product.select_one('.pricing-item-price').text.strip() #all the content of tag with price
16            product_url = product.select_one('.productPageLink').attrs['href'] #url of product page
17            try:
18                rating = product.select_one('.abt-reviews-rating .sr-only').text.strip()
19                review = product.select_one('.abt-reviews-num-reviews').text.strip()
20            except AttributeError:
21                rating = review = None
22                writer.writerow([url, category, title, price, product_url, rating, review])
23        nextbutton = "Next »" #scraping URL of next page
24        links = doc.find_all('a', string=nextbutton) #list of urls with only one element
25        if links:
26            url = links[0].attrs['href']
27            print('nxt page to scrape: ', url)
28            scrape_prods_to_csv(url, category) #simple recursion. Because why not

```

This procedure gathered all needed attributes of products (title with model ID, price, number of reviews and average rating) and the URL of next page. Except that function saved all the data into a file straight off.

Impurifying

Data already had natural impurities: duplicates, NA values, incorrect types, irrelevant data or contaminated data for some values, strings with words in prices (e.g. 'Your Price: \$799.00'), reviews and ratings ('5 out of 5 stars') where only numbers needed. But this is only 4 different types of impurities. So it was add two additional impurifications:

1. Removing the dot from prices to get outliers (float prices, in this case, could be treated as multiplied to 100) or misspells. Example 999.00\$ → 99900\$
2. Deleting titles. In case of this website, URL contains the same full title of the Product so the goal was to delete it and then recreate out of URL. Example of url:
<https://www.abt.com/Apple-iPad-9th-Generation-10.2-Inch-64GB-Wi-Fi-Space-Grey-2021-MK2K3LLA/p/168532.html>

Cleaning, Transformation and Uploading

During cleaning was removed artificial impurities, extracted correct prices, ratings and reviews, set to correct types. And then excessive columns were dropped,

During the transformation phase was added three additional calculations:

1. Price per GB
2. The average price of the same mode in different colours. Yes, they can differ.
3. Added price in CHF

Then prepared for merging data was uploaded to MariaDB.

Student C: Sasa Ljubisavljevic — www.cyberport.de

Extract

In the extraction phase, data was collected from the website <https://www.cyberport.de/apple-und-zubehoer.html> using the web scraper browser extension tool.



This tool allowed me to define selectors and create a structured scraping process. Templates were created to navigate the website's clear structure and target the desired Apple products. Subcategories were utilized to ensure precise item targeting and facilitate integration with other datasets.

The screenshot shows a web browser with a web scraping tool overlay. The main content is a product page for an 'Apple iPhone 14 Pro Max 128 GB Space Schwarz MQ9P3ZD/A'. The page includes a product image, a short description, a rating of 4 stars from 6 reviews, and a price of €1,269,00. Below the main content, there's a navigation bar with tabs like 'Elements', 'Console', 'Sources', etc., and a 'Web Scraper' tab which is currently selected. A large diagram below the navigation bar illustrates the data schema being scraped. The schema starts with a 'root' node, which branches into 'iphone_group' and 'iphone_model'. The 'iphone_group' node has a 'pagination' node. The 'iphone_model' node also has a 'pagination' node. Both 'pagination' nodes point to an 'element' node. This 'element' node then points to several data fields: 'model_id', 'price', 'rating', and 'ratings_number'. There are also three additional 'pagination' nodes branching off the 'element' node, each pointing to its own 'model_id', 'price', 'rating', and 'ratings_number' fields.

Impurifying and cleaning

To ensure data quality and remove unnecessary elements, impurities were addressed, and the data was cleaned.

Columns containing links for each element were dropped as they were not relevant to our analysis.

Rows with category values unrelated to our project were also removed. The data underwent thorough cleaning to standardize its format and eliminate unnecessary text and characters.

This included cleaning newly obtained columns such as storage and rating, ensuring a consistent format for analysis.

Transform

Following the extraction and cleaning stages, the data underwent transformation to make it suitable for analysis. Columns were modified, rearranged, and structured to facilitate further processing. Unnecessary text and characters were removed to create a refined dataset. Textual data was converted to numerical values to enable subsequent analysis and formatting.

Uploading

After the data was transformed and cleaned, it was uploaded to a MariaDB database. A table named "cyber_port" was created in the database to store the scraped data. This allowed for efficient storage and management of the dataset. Prior to uploading, a validation process was conducted to ensure the accuracy and integrity of the data. This step ensured that the uploaded data was ready for analysis and met the requirements of the project.

Merging

We could merge the data by following steps:

1. prepare the data from 3 source by same process:

- get the model 'number / storage / color' for the Iphone category.

brand	CH_title	CH_reviews	CH_rating	CH_price	category	storage	model	sub	prod_n	prod_n5
33 Apple	iPhone 13 - 512 GB, Red, 6.1", 12 MP, 5G	0	0.0	949.0	iPhone	512GB	13	Red	-1	-1
114 Apple	iPhone 11 - 64 GB, Black, 6.1", 12 MP, 4G	0	0.0	439.0	iPhone	64GB	11	Black	-1	-1
115 Apple	iPhone SE 3. Gen. - 64 GB, Midnight, 4.7" 12 M...	0	0.0	469.0	iPhone	64GB	SE 3	Midnight	-1	-1

b. create a column for the first 5 digits of the product number .

	brand	CH_title	CH_reviews	CH_rating	CH_price	category	storage	model	sub	prod_n	prod_n5
1	Apple	iPad Wi-Fi 2021 [10.2", 64 GB, Silver, MK2L3TY/A]	0	0.0	298.0	iPad	64GB	2021	10.2"	MK2L3TY/A	MK2L3
2	Apple	iPad Wi-Fi 2021 [10.2", 64 GB, Space Gray, MK2...	0	0.0	298.0	iPad	64GB	2021	10.2"	MK2K3TY/A	MK2K3
5	Apple	iPad Wi-Fi 2022 [10.9", 64 GB, Silver, MPQ03TY/A]	0	0.0	466.0	iPad	64GB	2022	10.9"	MPQ03TY/A	MPQ03

2. merge 3 data frame by the same process:

a. merge the CH and US data with the above 3 attributes for iPhone.

```
1 # merge the iphone data by storage, model, and color.
2 iphone_merged = pd.merge(df1,df2, left_on=['storage','model','sub'], right_on=['iphone_storage','iphone_model','iphone_color'])
3 iphone_merged
```

Python

	brand	CH_title	CH_reviews	CH_rating	CH_price	category	storage	model	sub	prod_n	prod_n5	US_title	model_id	iphone_storage	iphone_model	ip...
0	Apple	iPhone 12 - 64 GB, Black, 6.1", 12 MP, 5G	0	0.0	619.0	iPhone	64GB	12	Black	-1	-1	iPhone 12 Cellular Phone - MG...	MGEF3LL/A	64GB	12	...

b. merge the other category by short product number.

```
1 # merge the other product by first 5 model number.
2 other_merged = pd.merge(df1,df2, left_on=['prod_n5'], right_on=['model_id5'])
3 other_merged
```

Python

	brand	CH_title	CH_reviews	CH_rating	CH_price	category	storage	model	sub	prod_n	prod_n5	US_title
0	Apple	iPad Wi-Fi 2021 [10.2", 64 GB, Silver, MK2L3TY/A]	0	0.0	298.0	iPad	64GB	2021	10.2"	MK2L3TY/A	MK2L3	Apple iPad (9th Generation) 10.2-Inch 64GB Wi...

c. concat the 'iphone_merged' and 'other_merged' together.

3. subset the merged data with necessary columns. We are only interested in the price / reviews / ratings of 3 countries.

4. finally we create 2 columns to keep the 'CH_price_in_dollar' and "DE_price_in_dollar".

```
1 # take the exchange rate to exchange the price to dollar.
2 exchangerate = {'USD_CH':0.9046, 'EUR_USD': 1.0752}
3
4 # create two columns to keep the price in dollar.
5 final_merged['CH_price_in_dollar'] = (final_merged['CH_price']/exchangerate['USD_CH']).round(2)
6 final_merged['DE_price_in_dollar'] = (final_merged['DE_price'] * exchangerate['EUR_USD']).round(2)
7 final_merged
```

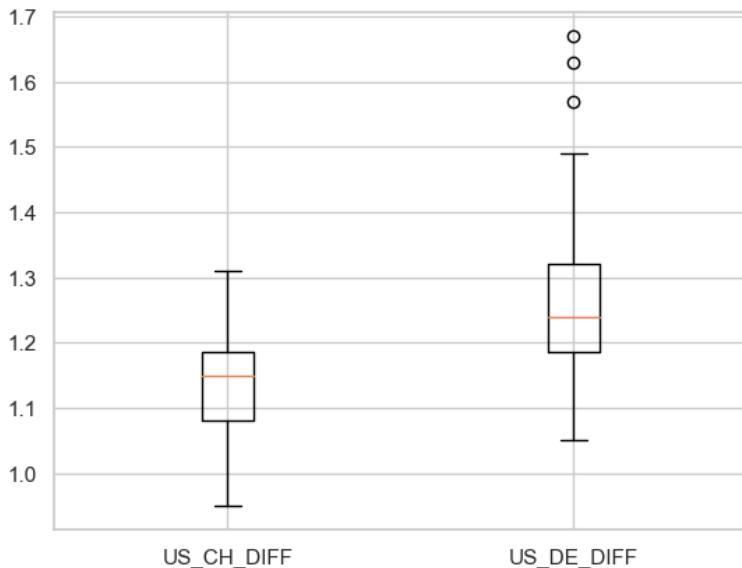
Analysis and answers to questions

Question 1: Which country has higher prices for Apple products, the USA or Switzerland? And by how much?

First, I use 'US_price' as the benchmark to normalize the price difference by following the code. and create 2 new columns to keep the results.

```
1 # normalize the price different by persentage.
2
3 final_merged['US_CH_DIFF'] = (final_merged['CH_price_in_dollar']/final_merged['US_price']).round(2)
4 final_merged['US_DE_DIFF'] = (final_merged['DE_price_in_dollar']/final_merged['US_price']).round(2)
```

Then we could plot the boxplot of those two columns to get our result.



results: the price in CH and DE both higher than the price in the US. On average, CH price is 1.15 times higher than price in the US; DE price is 1.25 times higher than the price in US.

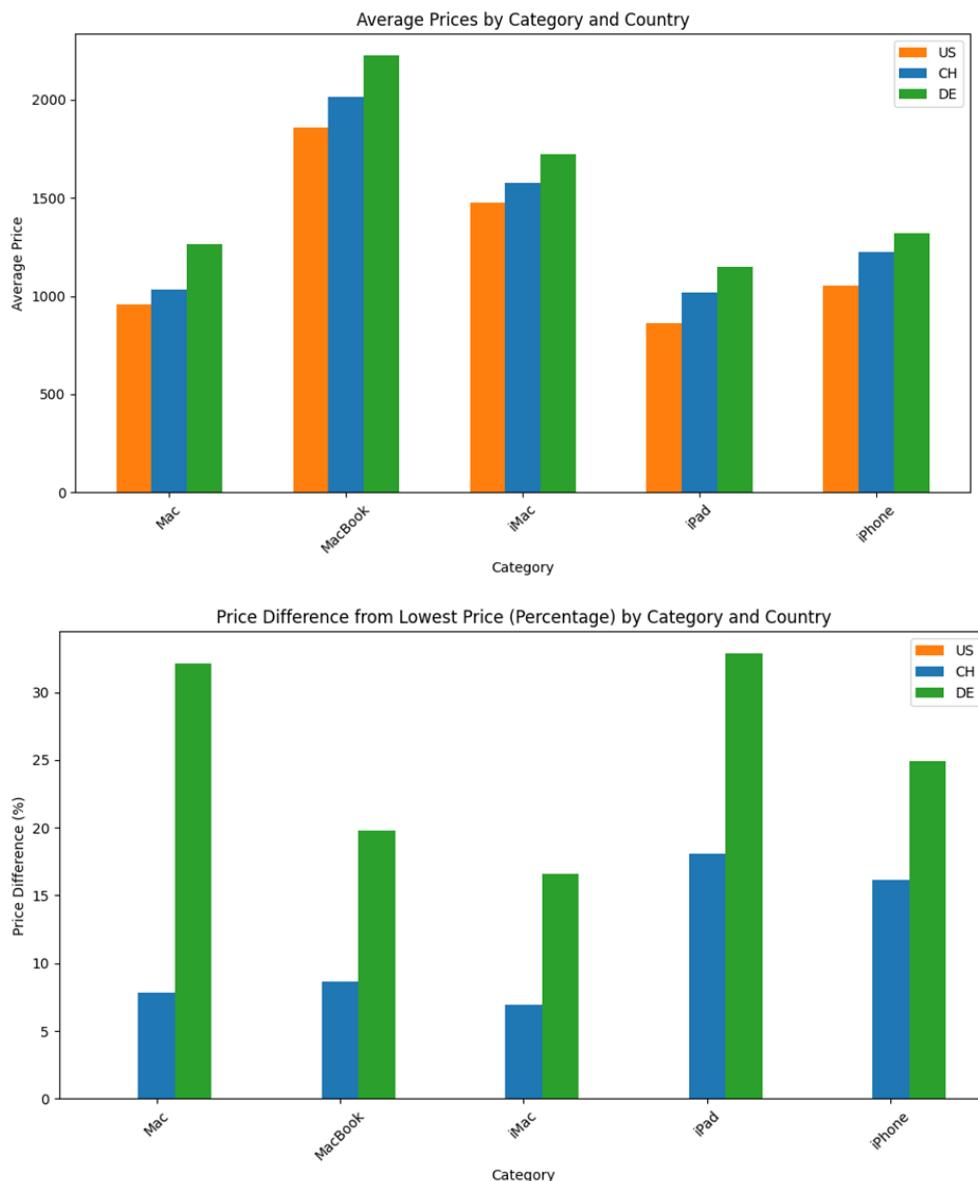
Question 2: Between Apple products categories, which has the biggest price difference?

In order to answer the second question, we approached in the following way: we grouped all articles from the final dataset and calculated their average value for each country. Then, in relation to the cheapest price from each category, we calculated the difference. Then we calculated their differences in relative values.

Category	Country	Average Price	Price Difference (Abs)	Price Difference (%)
13	Mac	US	956.200000	0.000000
12	Mac	CH	1030.954000	7.817821
14	Mac	DE	1263.358000	32.122778
7	MacBook	US	1856.200000	0.000000
6	MacBook	CH	2016.272800	8.623683
8	MacBook	DE	2223.682400	19.797565
10	iMac	US	1476.777778	0.000000
9	iMac	CH	1579.087778	6.927921
11	iMac	DE	1721.633333	16.580393
4	iPad	US	863.890364	0.000000
3	iPad	CH	1020.321091	18.107706
5	iPad	DE	1147.626545	32.844004
1	iPhone	US	1055.596061	0.000000
0	iPhone	CH	1225.871515	16.130740
2	iPhone	DE	1318.958939	24.949210

When comparing the average prices of Apple products across countries, the US generally offers the lowest prices for all categories. Switzerland tends to have moderately higher prices, while Germany consistently has the highest prices.

We then visualized the obtained results, which gives us a clear insight into the price differences between individual groups of items.



Results: The "Mac" and "iPad" category has the biggest price difference among Apple products. The price difference indicates the percentage difference in prices compared to the average price in the United States (US). The price difference value for the "Mac" and "iPad" category in Germany (DE) around 32%, which is the highest among all the categories listed.

Question 3: Is there a correlation or dependency between the popularity (in terms of Rating and number of reviews) and price differences of Apple products in the US and Germany?

We found some evidence that the difference in prices is dependent on Rating and number of reviews. P-value for both of them is less than 0.05. So within our scope with given shops, we found that:

- the larger rating the larger the relative difference in price. With every additional score the difference is growing on 13%.

- the larger the number of reviews the smaller relative difference in price we have. For one additional review we have drop in difference in prices for 0,18%

```

    OLS Regression Results
=====
Dep. Variable:      diff_price   R-squared:          0.175
Model:              OLS         Adj. R-squared:       0.154
Method:             Least Squares F-statistic:        8.352
Date:               Fri, 26 May 2023 Prob (F-statistic): 0.000512
Time:                23:45:24   Log-Likelihood:     72.140
No. Observations:    82         AIC:                 -138.3
Df Residuals:        79         BIC:                 -131.1
Df Model:            2
Covariance Type:    nonrobust
=====
      coef    std err      t      P>|t|      [ 0.025   0.975 ]
-----
const      -0.3115    0.299    -1.040    0.301    -0.908    0.285
all_reviews -0.0018    0.000    -3.626    0.001    -0.003    -0.001
avg_rating   0.1307    0.061    2.139    0.036     0.009    0.252
=====
Omnibus:           4.493   Durbin-Watson:        0.744
Prob(Omnibus):      0.106   Jarque-Bera (JB):   3.748
Skew:                0.420   Prob(JB):            0.154
Kurtosis:            2.376   Cond. No.          889.
=====
```

Limitations

We know that only three shops have been analysed and it is quite hard to draw conclusions about the whole country. Moreover, there was not so big number of reviews for chosen products in the shops analysed, so it is unlikely that the conclusion and answer to the third question can be extrapolated to the whole countries and shops. But the approach itself and the code created can be applied almost without changes to other shops and datasets, and conclusions can be drawn from these results.

Tools used

- Chrome and Developer tools in the Browser
- For scraping was used Jupyter Notebook. Main packages use:
 - Pandas
 - BeautifulSoup
 - Selenium
- Web scraper browser extension tool
- VM with Ubuntu
- MariaDB

Lesson learned / self-reflection

We learned, or rather confirmed the assumption, that it is most expensive to buy Apple products in Germany and cheapest in the USA. Found the category with the largest difference in prices. We found unexpected relationship between reviews and price differences, of course this data

needs to be double-checked. But it did not coincide with our hypothesis, and that in itself is interesting.

The next step in the development of the project could be scraping additional sites with more reviews and possibly in new countries. This would allow us to generalise the conclusions.

We have learnt how the basic tools for scrapping and data processing work, faced many challenges and solved many problems.

We have also seen the importance of starting work as soon as possible. And that even the seemingly easiest tasks can take an enormous amount of time, and one should be prepared for this. Additional study of supporting tools such as video editing applications and file editors can be very useful.

Overall, we found the project useful and instructive. Thanks for it!