

Generation

香 HONG KONG 港

Team Green

Jacky - Winson - Jeff



Data Analysis Tools

1. Market researchers more understand competition
2. Client's requirements and objectives for market
3. Identify the key data points (product price, promotion, update date etc..)
4. Determine the frequency of data collection and update schedule



Benefits (What's it bring?)

1. Comprehensive Market Insights
2. Real-time Data Collection
3. Competitive Advantage
4. Customization and Flexibility
5. Actionable Insights
6. Scalability and Future Growth



Brand



Wellcome 惠康

Wellcome is Hong Kong's the longest established supermarket chain. Wellcome's network of over 280 stores in Hong Kong.



PARKnSHOP 百佳

PARKnSHOP Supermarket is part of A.S. Watson Group, the world's largest international health and beauty retailer.



Ztore 士多

Ztore are an online supermarket with convenient delivery service, selling quality groceries and home essentials.

Scraping Tool

Wellcome 惠康

web crawler

提交意見 | 客戶幫助 | 註冊 | 登入 | EN | 中 | Powered by YUV

惠康網店 wellcome.com.hk

搜尋商品

買滿\$500即享免費送貨

添加收貨地址

分類

關注我們

門店位置 | 關於惠康

首頁 > 飲品

飲品

排序 | 推薦

多重優惠

經典口味 略爽滋味

Coca-Cola

\$38.00

可口可樂汽水罐裝 8 X 330ML

加入購物車

指定品牌送贈品

維他 低糖氣泡菊花茶 6 X 310ML

\$30.00

加入購物車

2件\$50

明治頂級鮮牛奶 946ML

\$33.00

加入購物車

多重優惠

無糖 無卡路里 略爽無比

Coca-Cola

\$38.00

可口可樂無糖汽水罐裝 8 X 330ML

加入購物車

指定分類送贈品

保利保鮮裝純牛奶 3 X 1LT

\$60.00 \$87.90

Paul's Milk

加入購物車

品牌

搜尋品牌

137 DEGREES

7 UP FREE 七喜輕怡

7 UP 七喜

ABBOTT 雅培

ACQUA PANNA 巴娜

展開

```
1 import pandas as pd
2 from selenium import webdriver
3 from selenium.webdriver.common.by import By
4 from selenium.webdriver.chrome.options import Options
5 import time
6 import datetime
```

Import the library which required
in web crawler

```
8 driver = webdriver.Chrome()
9 current_page= 1
10 url = "https://www.wellcome.com.hk/zh-hant/category/100002/" + str(current_page) + ".html"
11
12 try:
13     driver.get(url)
14 except:
15     print('url not found')
16     exit
17
18 time.sleep(2)
19
```

Use selenium to get the Wellcome
beverage website

Return error if the website is unavailable.



Create empty list and
dictionary for further use.

```
20   remark=[]
21   name = []
22   selling_price=[]
23   original_price=[]
24   link=[]
25   selling_price_base=[]

26
27   df_dict={
28       "name" : name,
29       "selling_price" : selling_price,
30       "original_price" : original_price,
31       "remark" : remark,
32       "link" : link
33   }
```

Find the last page
element in the website,
convert it to a integer
and store in a variable.

```
36   total_page = driver.find_element(By.XPATH, '//a[@class="last cursor num-box"]')
37   total_page = int(total_page.text)
38   #print('total page : ', total_page)
39
40
```



```
41 while total_page >= current_page:  
42     url = "https://www.wellcome.com.hk/zh-hant/category/100002/" + str(current_page) + ".html"  
43     try:  
44         driver.get(url)  
45         products = driver.find_elements(By.XPATH, '//a[@class="a-link router-link ware-wrapper"]')  
46         for product in products:  
47             # Process each product
```

Use while loop to access the first url(Page 1). Use XPATH to find the html class element which contain the whole product independently. Use for loop to obtain the data from the path.

The web crawling process and method will show on the next page.

```
        current_page = current_page+1  
    except:  
        print('url not found')  
        exit
```

At the end of the while loop, when all of the current page data is obtain, selenium will get in the next page



```
try:  
    remark.append(product.find_elements(By.XPATH, './div[contains(@class, "pro tag")][0].text)  
except:  
    remark.append(None)  
  
try:  
    name.append(product.find_elements(By.XPATH, './div[contains(@class, "name")][0].text)  
except:  
    name.append(None)  
  
try:  
    price_base = product.find_elements(By.XPATH, './div[contains(@class, "price")][0].text.split("\n")[0]  
    small_base = product.find_elements(By.XPATH, './div[contains(@class, "price")][0].text.split("\n")[1]  
    total=price_base+small_base  
    selling_price.append(total)  
  
except Exception as error:  
    selling_price.append(None)  
try:  
    original_price.append(product.find_elements(By.XPATH, './span[contains(@class, "line-price")][0].text)  
except:  
    original_price.append(None)  
  
product_link =driver.find_elements(By.XPATH, './a[contains(@class, "a-link router-link ware-wrapper")]')  
for product in product_link:  
    link.append(product.get_attribute('href'))  
  
current_page = current_page+1
```

In for loop, use XPATH to find the different part of data in the website, store that data with append in correspond list. If the data is empty or not found, the correspond list would append None as value.



```
87     df = pd.DataFrame(df_dict)
88
89     if __name__ == '__main__':
90         current_time_str = datetime.datetime.now().strftime("%Y_%m_%d")
91         print(df)
92         df.to_csv(current_time_str + "_wellcome.csv", encoding="UTF-8")
```

In for loop, use XPATH to find the different part of data in the website, store that data with append in correspond list. If the data is empty or not found, the correspond list would append None as value.

Ztore 士多 web crawler

賈滿 \$499 免費送貨上門

下載手機應用程式 | Eng | 新手指南 | Zmile Club及禮遇 | 士多生活誌

 搜尋產品名稱、類別、品牌或標籤

(健康店) (美妝店) 葱味蒜香撈麵 \$10/4包 吉百利朱古力任揀兩件半價 買米送米 買2送1 愛回家 x 一心形肉乾
轉季凍一凍 即食糖水\$35/2 KF94 \$45/60個

香港 登入/註冊

士多就係筍 終極感謝祭 店長推介 香港品牌 超市 保健醫藥 美容護理 母嬰

飲品

| | | | |
|---|---|--|---|
| 原箱飲品推介 172 即飲茶 - 紙包及罐裝、即飲茶 - 罐裝、汽水、水、果汁、椰子水... | 店長精選飲品 46 新登場飲品、期間限定飲品 | 即飲 - 茶類飲品、檸檬茶 192 檸茶 - 無糖茶、紙包及罐裝 - 無糖茶、罐裝 - 水果茶、甜茶、紙... | 即飲 - 咖啡、奶茶 19 咖啡、奶茶 |
| 有汽飲品 55 傳統汽水、特式汽水、梳打水、湯力水、奶類碳酸飲料、薑啤 | 水 48 礦泉水、礦物質水、蒸餾水、有氣水、水機套裝 | 果汁、椰子水 63 日韓產地名物果汁、瘦身養顏果汁、椰子水、蔬果汁、果汁-蘋... | 養生保健、花果茶 118 涼茶、草本飲品、養生保健飲品、米水、薑茶、黑糖粒粒、... |
| 植物奶、豆奶及豆浆 62 無糖植物奶、特別口味植物奶、杏仁奶、核桃奶、米奶、燕麥奶、... | 牛奶、淡奶、乳酪、奶粉 59 保鮮裝奶-全脂、保鮮裝奶-半脫/低脂、保鮮裝奶-全脫、有味奶... | 即沖飲品 90 即沖咖啡、掛耳式咖啡、咖啡粉囊、咖啡器具&咖啡機、咖啡粉... | 運動、能量飲品 29 能量飲品、葡萄糖、電解質飲品 |

澳洲 subuna 星空蘋果+ 獨家禮盒裝
研究證實

獨家禮盒裝
研究證實

코코베비 COCOBEBE 嬰兒濕紙巾-原箱
獨家禮盒裝
研究證實
共1000張!

Import Libraries

```
1 import pandas as pd
2 from selenium import webdriver
3 from selenium.webdriver.common.by import By
4 from selenium.webdriver.chrome.options import Options
5 import time
6 import datetime
```

Def Crawler

```
8 def crawler(url):
9     driver = webdriver.Chrome()
10    driver.get(url)
11    time.sleep(3)
```



Click the show all button and scroll down the web

```
25     try:
26         send_button = driver.find_element(By.XPATH, '//*[@id="BaseLayout"]/div/div[4]/span/div/span')
27         send_button.click()
28     except:
29         pass
30
31     time.sleep(3)
32
33     while len(data) < number:
34         data = driver.find_elements(By.XPATH, '//div[contains(@class, "jsx-2940394913 jsx-676457135 ProductItem windowing-layout")]')
35         driver.execute_script("window.scrollTo(0, document.body.scrollHeight);")
36         time.sleep(3)
```

Get Item Data

```
38     for index, item in enumerate(data, start=1):
39         try:
40             # item.find_element(By.CLASS_NAME, 'brand-product-name') found something
41             title_list.append(item.find_element(By.CLASS_NAME, 'name_bbd').text)
42         except:
43             # item.find_element(By.CLASS_NAME, 'brand-product-name') is not found
44             title_list.append(None)
45         try:
46             packsize_list.append(item.find_element(By.CLASS_NAME, 'packsize').text)
47         except:
48             packsize_list.append(None)
49         try:
50             OP_list.append(item.find_element(By.CLASS_NAME, 'price-container').text.split('\n')[0].split('$')[1])
51         except:
52             OP_list.append(None)
53         try:
54             DP_list.append(item.find_element(By.CLASS_NAME, 'price-container').text.split('\n')[-2].split('$')[1])
55         except:
56             DP_list.append(None)
```



+\$9.9換特製鹹柑桔 >

+\$25換柔感牙刷送姜... >

OOHA – 檸檬蜂蜜味汽水

★★★★★ (7)

\$15 \$3.75/罐

加入購物車



Get Promotion Data

[+\\$69.9換綜合維他命... >](#)

[+\\$25換柔感牙刷送姜... >](#)

[+\\$25換柔感牙刷送姜... >](#)

[+\\$69.9換綜合維他命... >](#)

[+\\$25換柔感牙刷送姜... >](#)

[+\\$69.9換綜合維他命... >](#)

```
58     try:
59         promotions = item.find_elements(By.XPATH, f"//*[@id='BaseLayout']/div/div[4]/div[2]/div/div/div[{index}]/div[3]/div/div")
60         promotions = promotions[0].text.split("\n")
61
62         try:
63             Promotion1_list.append(promotions[0])
64         except:
65             Promotion1_list.append(None)
66
67         try:
68             Promotion2_list.append(promotions[1])
69         except:
70             Promotion2_list.append(None)
71     except:
72         Promotion1_list.append(None)
73         Promotion2_list.append(None)
74
75     try:
76         link_list.append(item.find_element(By.XPATH, f'//*[@id="BaseLayout"]/div/div[4]/div[2]/div/div/div[{index}]/a').get_attribute("href"))
77     except:
78         link_list.append(None)
```

Dataframe Design

```
81     df = pd.DataFrame(  
82         {  
83             'Product_Title' : title_list,  
84             'Package_Size' : packsize_list,  
85             'Original_Price' : OP_list,  
86             'Discount_Price' : DP_list,  
87             'Promotion_1' : Promotion1_list,  
88             'Promotion_2' : Promotion2_list,  
89             'Link' : link_list  
90         }  
91     )  
92  
93     df['Discount_Price'] = df['Discount_Price'].replace(',', '', regex=True).astype(float)  
94     df['Original_Price'] = df['Original_Price'].replace(',', '', regex=True).astype(float)  
95  
96     driver.quit()  
97  
98     return df
```

Dictionary URL and Categories

```

100 zstore_drinks_url = {
101     "Beverage Case Deals" : "https://www.zstore.com/tc/category/all/beverage/case-offers",
102     "Recommended Drinks" : "https://www.zstore.com/tc/category/all/beverage/recommended-drinks",
103     "RTD - Tea & Lemon Tea" : "https://www.zstore.com/tc/category/all/beverage/rtd-tea-lemon-tea",
104     "RTD - Coffee & Milk Tea" : "https://www.zstore.com/tc/category/all/beverage/rtd-coffee-milk-tea",
105     "Carbonated Drinks" : "https://www.zstore.com/tc/category/all/beverage/carbonated-beverage",
106     "Water" : "https://www.zstore.com/tc/category/all/beverage/water",
107     "Juice and Coconut Water" : "https://www.zstore.com/tc/category/all/beverage/juice-and-energy-drink",
108     "Healthy Drinks" : "https://www.zstore.com/tc/category/all/beverage/healthy-drinks",
109     "Plant Based & Soy Milk" : "https://www.zstore.com/tc/category/all/beverage/soy-plant-based-milk",
110     "Long Life Milk, Yogurt & Milk Powder" : "https://www.zstore.com/tc/category/all/beverage/long-life-milk-soy",
111     "Hot Beverage" : "https://www.zstore.com/tc/category/all/beverage/coffee-tea",
112     "Energy & Sports Drink" : "https://www.zstore.com/tc/category/all/beverage/energy-drink"
113 }

```

飲品

| | | | |
|--|--|--|---|
| 原箱飲品推介 169 即飲茶 – 紙包及罐裝、即飲茶 – 樋裝、汽水、水、果汁、椰子水... | 店長精選飲品 45 新登場飲品、期間限定飲品 | 即飲 – 茶類飲品、檸檬茶 189 樋裝 – 無糖茶、紙包及罐裝 – 無糖茶、檳榔茶 – 水果茶、甜茶、紙... | 即飲 – 咖啡、奶茶 19 咖啡、奶茶 |
| 有汽飲品 48 傳統汽水、特式汽水、梳打水、湯力水、奶類碳酸飲料、薑啤 | 水 48 礦泉水、礦物質水、蒸餾水、有氣水、水機套餐 | 果汁、椰子水 61 日韓產地名物果汁、瘦身養顏果汁、椰子水、蔬果汁、果汁-蘋... | 養生保健、花果茶 115 涼茶、草本飲品、養生保健飲品、米水、薑棗、薑茶、黑糖粒粒、... |
| 植物奶、豆奶及豆浆 61 無糖植物奶、特別口味植物奶、杏仁奶、核桃奶、米奶、燕麥奶、... | 牛奶、淡奶、乳酪、奶粉 59 保鮮裝奶 – 全脂、保鮮裝奶 – 半脫/低脂、保鮮裝奶 – 全脫、有味奶... | 即沖飲品 88 即沖咖啡、掛耳式咖啡、咖啡粉囊、咖啡器具&咖啡機、咖啡粉... | 運動、能量飲品 27 能量飲品、葡萄糖、電解質飲品 |



Print csv

```
117     # Create current date (datetime obj), and change it into string.  
118     current_time_str = datetime.datetime.now().strftime("%Y_%m_%d")  
119  
120     for k, v in ztore_drinks_url.items():  
121         df = crawler(v)  
122         # Use assign to create a new column "date" with the same value (current_datetime)  
123         df = df.assign(date = current_time_str)  
124  
125         # Create a new column "retailer" with retailer name "Ztore"  
126         df = df.assign(retailer = "ztore")  
127  
128         # Create a new column "category" with the pns_drinks_url key  
129         df = df.assign(category = k)  
130  
131         df.to_csv(k + "_" + current_time_str + "_ztore.csv", encoding="UTF-8")
```

PARKnSHOP 百佳

web crawler



There are 11 categories of beverages

飲品、即沖飲品

[查看所有](#)

| | | | | | | | |
|---|---------|---|----------|--|------------|---|------------|
|  | 酒精飲品 |  | 水 |  | 汽水 |  | 即飲茶類、咖啡、奶茶 |
|  | 奶類、乳酪飲品 |  | 植物奶、大豆飲品 |  | 咖啡、沖調飲品、熱飲 |  | 果汁、椰子水 |
|  | 運動及能量飲品 |  | 草本及健康飲品 |  | 原箱飲品 | | |

PARKnSHOP 百佳 web crawler

The screenshot shows the PARKnSHOP (PNS) website interface. At the top, there is a navigation bar with links for '易賞鑑VIP會員獎賞' (Easy賞鑑VIP Member Reward), '店舖一覽' (Shop List), and 'English'. The main menu includes categories like '食品及飲品' (Food & Beverage), '母嬰' (Mother & Baby), '個人護理、健康' (Personal Care, Health), '家居生活' (Home Living), '寵物食品及護理' (Pet Food & Care), '獨家產品' (Exclusive Products), '護膚美妝' (Skincare & Beauty), and '屈臣氏旗艦店' (Watson's Flagship Store). A search bar is located at the top center.

The current page displays the '汽水' (Soda) category under the 'Food & Beverage' section. The breadcrumb navigation shows: 主頁 > 食品及飲品 > 飲品、即沖飲品 > 汽水.

The page features a product filter section with tabs for '分類' (Category), '品牌' (Brand), '價格' (Price), '推廣優惠' (Promotion), '原產地' (Origin), and '進階選項' (Advanced Options). Below this are four product categories: '全部' (All) featuring Coca-Cola, '汽水' (Soda) featuring Coca-Cola, '運動、能量飲品' (Sports, Energy Drinks) featuring Red Bull, and '梳打水、漏力水' (Soda, Lemonade) featuring Watson's.

The main content area shows a grid of 160 products. The first few items are highlighted with 'BEST' awards and promotional offers:

- 屈臣氏 卡路里0卡路里 Watson's ENERGY DRINK 组合優惠 (Offer)
- 屈臣氏 蘋果味蘇打水 送贈品 (Offer)
- 屈臣氏 低卡路里 Watson's ENERGY DRINK 组合優惠 (Offer)
- 屈臣氏 屈臣氏青檸味蘇打水 Watson's ENERGY DRINK 组合優惠 (Offer)
- 屈臣氏 屈臣氏湯力水 Watson's TONIC WATER 送贈品 (Offer)
- 屈臣氏 屈臣氏蘇打水330 Watson's SODA 送贈品 (Offer)

Page layout

- Contains many products in the category
- When scrolldown, page will shows more products, until no more product



```
import pandas as pd
from selenium import webdriver # use webdriver
from selenium.webdriver.common.by import By # different methods of locating data
from selenium.webdriver.chrome.options import Options # options for selenium driver
import time
import datetime
```

Import the libraries which required

```
pns_drinks_url = {
    "Alcoholic Beverages" : "https://www.pns.hk/zh-hk/drinks/alcohol",
    "Water" : "https://www.pns.hk/zh-hk/drinks/water",
    "Carbonated Drink" : "https://www.pns.hk/zh-hk/drinks/carbonated-drink",
    "Instant Tea & Coffee & Milk Tea" : "https://www.pns.hk/zh-hk/drinks/tea-coffee-milk-tea",
    "Milk & Yogurt" : "https://www.pns.hk/zh-hk/drinks/milk-yogurt",
    "Plant Based & Soy Milk" : "https://www.pns.hk/zh-hk/drinks/plant-based-soy-milk",
    "Coffee, Hot Drink & Mix Powder" : "https://www.pns.hk/zh-hk/drinks/coffee-hot-drink-mix-powder",
    "Juice & Coconut Water" : "https://www.pns.hk/zh-hk/drinks/juice-coconut-water",
    "Energy Drink" : "https://www.pns.hk/zh-hk/drinks/energy-drink",
    "Herbal & Healthy Drink" : "https://www.pns.hk/zh-hk/drinks/herbal-healthy-drink",
    "Case Offer" : "https://www.pns.hk/zh-hk/drinks/case-offer"
}
```

Create dictionary contains categories and URLs

```

import pandas as pd
from selenium import webdriver # use webdriver
from selenium.webdriver.common.by import By # different methods of locating data
from selenium.webdriver.chrome.options import Options # options for selenium driver
import time
import datetime

> def crawler(url): ...

pns_drinks_url = {
    "Alcoholic Beverages" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E9%85%92%E7%B2%BE%E9%A3%81",
    "Water" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E6%80%84/c/04010100",
    "Carbonated Drink" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E6%80%84/c/04010100",
    "Instant Tea & Coffee & Milk Tea" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E5%8D%B3%E9%A3%B2%E8%8C%8A",
    "Milk & Yogurt" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E5%A5%86%E9%A1%9E%E3%80%80",
    "Plant Based & Soy Milk" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E6%A4%8D%7%89%A9%E5%A5%80%80",
    "Coffee, Hot Drink & Mix Powder" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E5%92%96%E5%95%A1%E3%80%80",
    "Juice & Coconut Water" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E6%9E%9C%E6%81%81%E3%80%80",
    "Energy Drink" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E9%81%88%E5%8B%95%E5%8F%80%80",
    "Herbal & Healthy Drink" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E8%8D%89%E6%9C%AC%E5%8F%80%80",
    "Case Offer" : "https://www.pns.hk/zh-hk/%E9%A3%9F%E5%93%81%E5%8F%8A%E9%A3%B2%E5%93%81/%E5%8E%9F%E7%AE%B1%E9%A3%81"
}

# create current date (datetime obj), and change it into string
current_time_str = datetime.datetime.now().strftime("%Y_%m_%d")

for k, v in pns_drinks_url.items():
    df = crawler(v)
    # use assign to create a new column "date" with the same value (current_datetime)
    df = df.assign(date = current_time_str)

    # create a new column "retailer" with retailer name "parknshop"
    df = df.assign(retailer = "parknshop")

    # create a new column "category" with the pns_drinks_url key
    df = df.assign(category = k)

    # save as csv file
    df.to_csv("pns_data/" + k + "_" + current_time_str + "_pns.csv", encoding="UTF-8")

```



Dictionary contains categories and URLs

Loop through all URLs

Use crawler function to get data
input :url
return :dataframe

Insert today's date and category into the columns

Save as csv files named with category and date



```
def crawler(url):

    driver = webdriver.Chrome() # create instance of selenium driver

    try:
        driver.get(url)
    except:
        print("url not found")
        exit

    time.sleep(2)

    remark = []
    name = []
    capacity = []
    selling_price = []
    original_price = []
    more_info = []
    link = []

    result_dict = {
        "remark" : remark,
        "name" : name,
        "capacity" : capacity,
        "selling_price" : selling_price,
        "original_price" : original_price,
        "more_info" : more_info,
        "link" : link
    }
```

Define crawler function

Try to the web page
exit when url not found

Create empty lists and dictionary
for the date we will get

```

try:
    total_no_of_product = driver.find_element(By.XPATH, "/html/body/app-root/cx-storefront/main/cx-page-layout/cx-page-slot//pns-product-list/div[1]/div[1]")
    total_no_of_product = int(total_no_of_product.text.split(" ")[0])
    print(total_no_of_product)

    products = driver.find_elements(By.XPATH, "//html/body/app-root/cx-storefront/main/cx-page-layout/cx-page-slot//pns-product-list/div[2]/div/pns-product-tile")
    print(len(products))

    while len(products) < total_no_of_product:
        driver.execute_script("window.scrollTo(0, document.body.scrollHeight)")
        time.sleep(2)
        products = driver.find_elements(By.XPATH, "//html/body/app-root/cx-storefront/main/cx-page-layout/cx-page-slot//pns-product-list/div[2]/div/pns-product-tile")
        print(len(products))

```



篩選條件:

160 件貨品

熱門 最新 評分 優惠 價錢 ↑↓


屈臣氏 蘇打水4罐裝(隨機發貨)
330MLX4


送贈品


屈臣氏 屈臣氏青檸味蘇打水 330MLX4


屈臣氏 屈臣氏湯力水 330MLX4


屈臣氏 屈臣氏蘇打水330毫升(原味/美希托味/四罐混合裝)(隨機發貨)
4X330ML

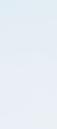

樂天 忌廉溝鮮奶(奶類碳酸飲品)
250MLx6


可口可樂 可口可樂8罐裝
330MLX8


可口可樂 無糖可口可樂汽水
330MLX12


樂天 忌廉溝鮮奶
500ML


可口可樂 可口可樂8罐裝
330MLX8


可口可樂 無糖可口可樂汽水
330MLX12

Find the number of products
scrolldown web page to show more products

```

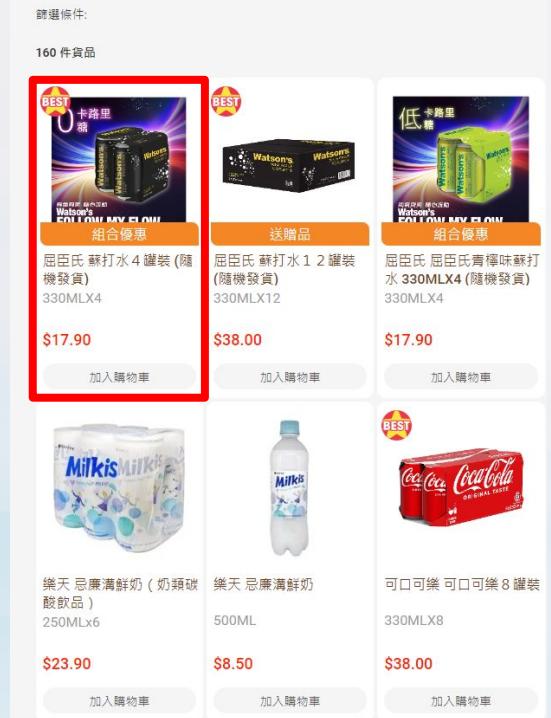
for product in products:
    product_info = product.text.split("\n")
    if len(product_info) == 4:
        remark.append(None)
        name.append(product_info[0])
        capacity.append(product_info[1])
        selling_price.append(product_info[2])
        original_price.append(None)
        more_info.append(product_info[-1])

    elif len(product_info) == 5:
        if "$" in product_info[-2] and "$" in product_info[-3]:
            remark.append(None)
            name.append(product_info[0])
            capacity.append(product_info[1])
            selling_price.append(product_info[2])
            original_price.append(product_info[3])
        else:
            remark.append(product_info[0])
            name.append(product_info[1])
            capacity.append(product_info[2])
            selling_price.append(product_info[3])
            original_price.append(None)

        more_info.append(product_info[-1])

    elif len(product_info) == 6:
        remark.append(product_info[0])
        name.append(product_info[1])
        capacity.append(product_info[2])
        selling_price.append(product_info[3])
        original_price.append(product_info[4])
        more_info.append(product_info[-1])

```



Get data and append to the correspond lists

```

product_link = driver.find_elements(By.XPATH, "/html/body/app-root/cx-storefront/main/cx-page-layout/cx-page-slot//pns-product-list/div[2]/div/pns-product-tile//div/div[1]/div[2]/a")
for product in product_link:
    link.append(product.get_attribute('href'))

except Exception as error:
    print(f'Error - {error} = {url}')

```



```
driver.quit()

df = pd.DataFrame(result_dict)

df['capacity'] = df['capacity'].str.upper()      # change capacity column to uppercase

df.index.name = 'sequence'          # the index column of df will be named 'sequence'

return df
```

```
# create current date (datetime obj), and change it into string
current_time_str = datetime.datetime.now().strftime("%Y_%m_%d")

for k, v in pns_drinks_url.items():
    df = crawler(v)
    # use assign to create a new column "date" with the same value (current_datetime)
    df = df.assign(date = current_time_str)

    # create a new column "retailer" with retailer name "parknshop"
    df = df.assign(retailer = "parknshop")

    # create a new column "category" with the pns_drinks_url key
    df = df.assign(category = k)

    # save as csv file
    df.to_csv("pns_data/" + k + "_" + current_time_str + "_pns.csv", encoding="UTF-8")
```

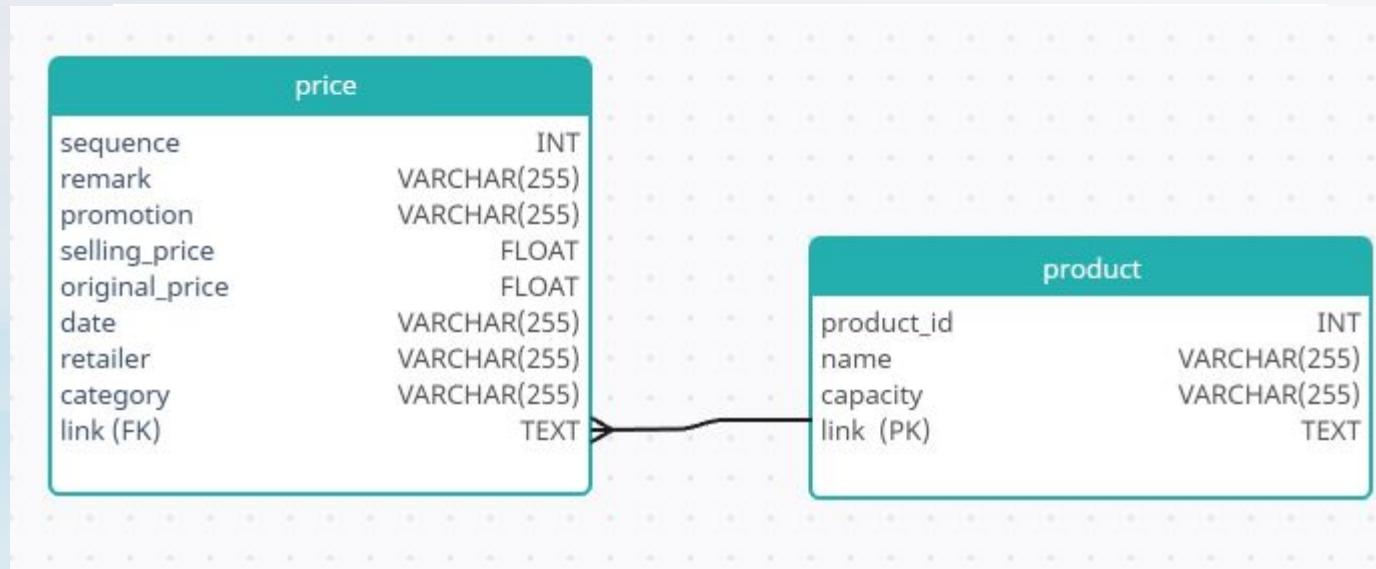
The function returns dataframe

Add current date, retailer name and category into the dataframe

Save as csv file

Database

Entity Relationship Diagram and Schemas



Create two tables
Decrease data redundancy

Create table and insert values

```
import psycopg2
import pandas as pd
import os
```

Import the libraries required

```
# Create a connection to the database
conn = psycopg2.connect(
    host = "...",
    dbname = "...",
    user = "...",
    password = "...",
    port = ...)
```

Connect to the database such as PostgreSQL

Create table

```
cur.execute(  
    """  
    DROP TABLE IF EXISTS product CASCADE;  
    CREATE TABLE product (  
        product_id INT,  
        name VARCHAR(255),  
        capacity VARCHAR(255),  
        link TEXT PRIMARY KEY  
    );  
  
    DROP TABLE IF EXISTS price CASCADE;  
    CREATE TABLE price (  
        sequence INT,  
        remark VARCHAR(255),  
        promotion VARCHAR(255),  
        selling_price FLOAT,  
        original_price FLOAT,  
        date VARCHAR(255),  
        retailer VARCHAR(255),  
        category VARCHAR(255),  
        link TEXT,  
        FOREIGN KEY (link) REFERENCES product(link)  
    );  
  
    """")
```

Use SQL query in python to create tables in database

Insert values

```
masterdf = pd.DataFrame()

for file in files:
    if file.endswith(".csv"):
        df = pd.read_csv(dataFolder + file)
        df.rename(columns={"Unnamed: 0":'sequence',
                           "Link" : "link",
                           "Product_Title" : "name",
                           "Package_Size" : "capacity",
                           "Original_Price" : "original_price",
                           "Discount_Price" : "selling_price",
                           "Promotion_1" : "remark",
                           "Promotion_2" : "promotion"
                         }, inplace=True)

masterdf = pd.concat((masterdf, df), ignore_index = True)
```

Get files from folder

Use Python logical conditions
to read only csv files

Insert values

```
# Specify the columns you want to insert
columns_to_insert_product = [
    "product_id",
    "name",
    "capacity",
    "link"
]

columns_to_insert_price = [
    "sequence",
    "remark",
    "promotion",
    "selling_price",
    "original_price",
    "date",
    "retailer",
    "category",
    "link"
]
```

Insert values

```
# Create a cursor object
cur = conn.cursor()

# Insert into the table

for index, row in masterdf.iterrows():
    # Construct the INSERT INTO query
    query = f"INSERT INTO product ({', '.join(columns_to_insert_product)}) VALUES ({', '.join(['%' + 's' * len(columns_to_insert_product))])} ON CONFLICT (link) DO NOTHING"

    # Prepare the values
    values = [row[column] for column in columns_to_insert_product]

    # Execute the INSERT INTO query
    cur.execute(query, values)

    query = f"INSERT INTO price ({', '.join(columns_to_insert_price)}) VALUES ({', '.join(['%' + 's' * len(columns_to_insert_price))})}"
    values = [row[column] for column in columns_to_insert_price]
    cur.execute(query, values)
```

```
# Commit the changes
conn.commit()

# Close the connection
conn.close()
```

Use SQL query in python to insert data into database

Price table schema

Query Query History

```
1 select * from price;
```

| price | |
|------------|--------------|
| product_id | INT |
| name | VARCHAR(255) |
| capacity | VARCHAR(255) |
| link (PK) | TEXT |

Data Output Messages Notifications

| | sequence | remark | promotion | selling_price | original_price | date | retailer | category | link |
|----|----------|-------------------------|-------------------------|------------------|------------------|-------------------------|-------------------------|-------------------------|-------|
| | integer | character varying (255) | character varying (255) | double precision | double precision | character varying (255) | character varying (255) | character varying (255) | text |
| 1 | 0 | 多重优惠 | [nul] | | 38 | NaN | 2023_11_17 | wellcome | [nul] |
| 2 | 1 | 指定品牌送赠品 | [nul] | | 30 | NaN | 2023_11_17 | wellcome | [nul] |
| 3 | 2 | 2件\$50 | [nul] | | 33 | NaN | 2023_11_17 | wellcome | [nul] |
| 4 | 3 | 多重优惠 | [nul] | | 38 | NaN | 2023_11_17 | wellcome | [nul] |
| 5 | 4 | 指定品牌送赠品 | [nul] | | 26 | NaN | 2023_11_17 | wellcome | [nul] |
| 6 | 5 | NaN | [nul] | | 13 | NaN | 2023_11_17 | wellcome | [nul] |
| 7 | 6 | 多重优惠 | [nul] | | 38 | NaN | 2023_11_17 | wellcome | [nul] |
| 8 | 7 | 多重优惠 | [nul] | | 18 | NaN | 2023_11_17 | wellcome | [nul] |
| 9 | 8 | NaN | [nul] | | 13 | NaN | 2023_11_17 | wellcome | [nul] |
| 10 | 9 | 指定分類送赠品 | [nul] | | 13 | NaN | 2023_11_17 | wellcome | [nul] |
| 11 | 10 | NaN | [nul] | | 12.9 | NaN | 2023_11_17 | wellcome | [nul] |
| 12 | 11 | 多重优惠 | [nul] | | 38 | NaN | 2023_11_17 | wellcome | [nul] |
| 13 | 12 | NaN | [nul] | | 13 | NaN | 2023_11_17 | wellcome | [nul] |
| 14 | 13 | NaN | [nul] | | 24 | 29 | 2023_11_17 | wellcome | [nul] |
| 15 | 14 | 指定品牌送赠品 | [nul] | | 26 | NaN | 2023_11_17 | wellcome | [nul] |
| 16 | 15 | 2件\$41 | [nul] | | 23 | NaN | 2023_11_17 | wellcome | [nul] |
| 17 | 16 | 2件\$32 | [nul] | | 20 | NaN | 2023_11_17 | wellcome | [nul] |

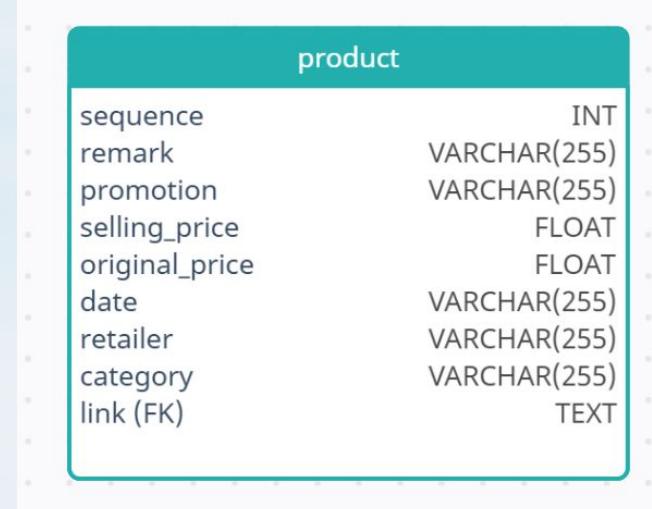
Product table schema

Query Query History

```
1 select * from product;
```

Data Output Messages Notifications

| | product_id | name | capacity | link |
|----|------------|-----------------------|----------|---|
| 1 | [null] | 可口可樂汽水罐裝 8 X 330ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E5%8F%AF%E5%AD%90 |
| 2 | [null] | 维他 低糖氣泡菊花茶 6 X 310ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 3 | [null] | 明治頂級鮮牛乳 946ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E6%98%8E%E6%9C%8B |
| 4 | [null] | 可口可樂無糖汽水罐裝 8 X 330ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E5%8F%AF%E5%AD%90 |
| 5 | [null] | 维他檸檬茶 9 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 6 | [null] | 维他冷泡無糖香片茶 6 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 7 | [null] | 可口可樂加系汽水罐裝 8 X 330ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E5%8F%AF%E5%AD%90 |
| 8 | [null] | 维他檸檬茶 6 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 9 | [null] | 维他冷泡無糖凍頂烏龍茶 6 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 10 | [null] | 北海道3.6牛乳迷你裝 200ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E5%8C%97%E6%9D%A1 |
| 11 | [null] | 维他蘋果綠茶 6 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 12 | [null] | 玉泉罐裝梳打水 8 X 330ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%8E%89%E6%9D%A1 |
| 13 | [null] | 维他冷泡無糖纖體音茶 6 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 14 | [null] | Meadows全脂鮮牛乳 1LT | [null] | https://www.wellcome.com.hk/zh-hant/p/Meadows%E5%85%A8 |
| 15 | [null] | 维他低糖檸檬茶 9 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 16 | [null] | 屈臣氏蒸餲水 4.5LT | [null] | https://www.wellcome.com.hk/zh-hant/p/%E5%B1%88%E8%99%9F |
| 17 | [null] | MEADOWS 純打水 6 X 320ML | [null] | https://www.wellcome.com.hk/zh-hant/p/MEADOWS%20%E6%9D%A1 |
| 18 | [null] | 维他高鈣朱古力奶 4 X 125ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |
| 19 | [null] | 维他冷泡無糖錫蘭檸檬茶 6 X 250ML | [null] | https://www.wellcome.com.hk/zh-hant/p/%E7%B6%AD%E4%9C%A1 |



Join two tables for analyzing data

Query Query History Scratch Pad

```
1 select * from product
2 join price
3 using (link);
```

Data Output Messages Notifications

| | link text | product_id integer | name character varying (2 | capacity character varying | sequence integer | remark character varyin | promotion character varyin | selling_price double precision | original_price double precisor | date character vari | retailer character varyi | category character varying (255) |
|----|-------------------|-----------------------|------------------------------|-------------------------------|---------------------|----------------------------|-------------------------------|-----------------------------------|-----------------------------------|------------------------|-----------------------------|-------------------------------------|
| 1 | https://www.pn... | 327901 | 葡萄适 葡萄适 S... | 750ML | 22 | \$24.9 / 2件 | [null] | 14.9 | NaN | 2023_11_... | parknshop | Energy Drink |
| 2 | https://www.pn... | 45321 | 好茶養生 桂花雪... | EACH | 49 | NaN | [null] | 15 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 3 | https://www.pn... | 45323 | 好茶養生 寧神安... | EACH | 50 | NaN | [null] | 12 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 4 | https://www.pn... | 45325 | 好茶養生 桂圓紅... | EACH | 51 | NaN | [null] | 12 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 5 | https://www.pn... | 45326 | 好茶養生 清熱降... | EACH | 52 | NaN | [null] | 12 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 6 | https://www.pn... | 45327 | 好茶養生 茉莉玫... | EACH | 53 | NaN | [null] | 12 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 7 | https://www.pn... | 45329 | 好茶養生 皇牌老... | EACH | 54 | NaN | [null] | 180 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 8 | https://www.pn... | 45330 | 好茶養生 黑糖薑... | EACH | 55 | NaN | [null] | 180 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 9 | https://www.pn... | 45335 | 好茶養生 枸杞阿... | EACH | 56 | NaN | [null] | 168 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 10 | https://www.pn... | 45337 | 好茶養生 黑棗十... | EACH | 57 | NaN | [null] | 168 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 11 | https://www.pn... | 45352 | 好茶養生 仙楂烏... | EACH | 58 | NaN | [null] | 50 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 12 | https://www.pn... | 45354 | 好茶養生 玫瑰四... | EACH | 59 | NaN | [null] | 50 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 13 | https://www.pn... | 45362 | 好茶養生 蜜糖菊... | EACH | 60 | NaN | [null] | 50 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 14 | https://www.pn... | 45381 | 好茶養生 蜂蜜檸... | EACH | 61 | NaN | [null] | 80 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |
| 15 | https://www.pn... | 48648 | 好茶養生 石斛花... | EACH | 62 | NaN | [null] | 15 | NaN | 2023_11_... | parknshop | Coffee, Hot Drink & Mix ... |

Analyzing data

Example : When the product release

```
import pandas as pd
import psycopg2
import matplotlib.pyplot as plt
from matplotlib.font_manager import FontProperties
```

```
# Establish a connection to the database
> conn = psycopg2.connect(...)

# SQL query
sql_query = """
    SELECT * FROM price
    JOIN product
    USING (link)
    where name like 'BOSS%';
"""

# Use pandas to execute the SQL query and store the result in a DataFrame
df = pd.read_sql_query(sql_query, conn)

# close the connection
conn.close()
```

Using matplotlib to create diagram

Analyzing data

example :to show when the product release

```
# Specify a font that supports Chinese characters
font = FontProperties(fname=r"font.ttf") # the path to your font file

df['date'] = pd.to_datetime(df['date'], format="%Y-%m-%d") # Ensure the date column is in datetime format

# Sort the DataFrame based on date
df.sort_values('date', inplace=True)

# Create a new column 'available' where each product is 1 if available and 0 if not
df['available'] = df['name'].notna().astype(int)

start_date = '2023-11-01'
end_date = '2023-11-18'

# Create a date range that includes all dates you're interested in
all_dates = pd.date_range(start=start_date, end=end_date)

# Loop over each product
for product in df['name'].unique():
    # Filter data for the current product
    product_data = df[df['name'] == product]

    # Reindex the DataFrame to include all dates and fill missing availability with 0
    product_data.set_index('date', inplace=True)
    product_data = product_data.reindex(all_dates, fill_value=0)

    # Reset the index so 'date' is a column again
    product_data.reset_index(inplace=True)
    product_data.rename(columns={'index': 'date'}, inplace=True)

    # Plot the data
    plt.plot(product_data['date'], product_data['available'], label=product)
```

using matplotlib to create diagram

Analyzing data

Example : When the product release

Using matplotlib to create diagram

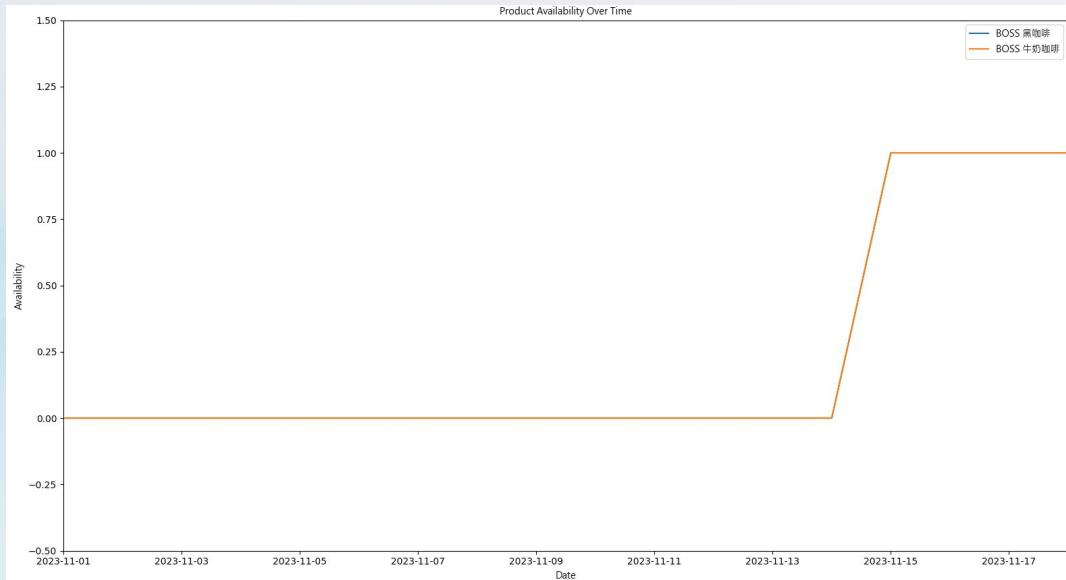
```
# Set the range of the x-axis
plt.xlim(pd.to_datetime(start_date, format="%Y-%m-%d"), pd.to_datetime(end_date, format="%Y-%m-%d"))

# Set the range of the y-axis
plt.ylim(-0.5, 1.5)

plt.xlabel('Date', fontproperties=font)
plt.ylabel('Availability', fontproperties=font)
plt.title('Product Availability Over Time', fontproperties=font)
plt.legend(prop=font)
plt.show()
```

Analyzing data

Example : When the product release output:



Result:
The brand BOSS 黑咖啡, 牛奶咖啡
released on 16 Nov 2023

Analyzing data

Example : When the product release

Query Query History

```
1 select category, count(*) as no_of_product from price
2 where date = '2023_11_17'
3 and retailer = 'parknshop'
4 group by category;
```

Data Output Messages Graph Visualiser Notifications

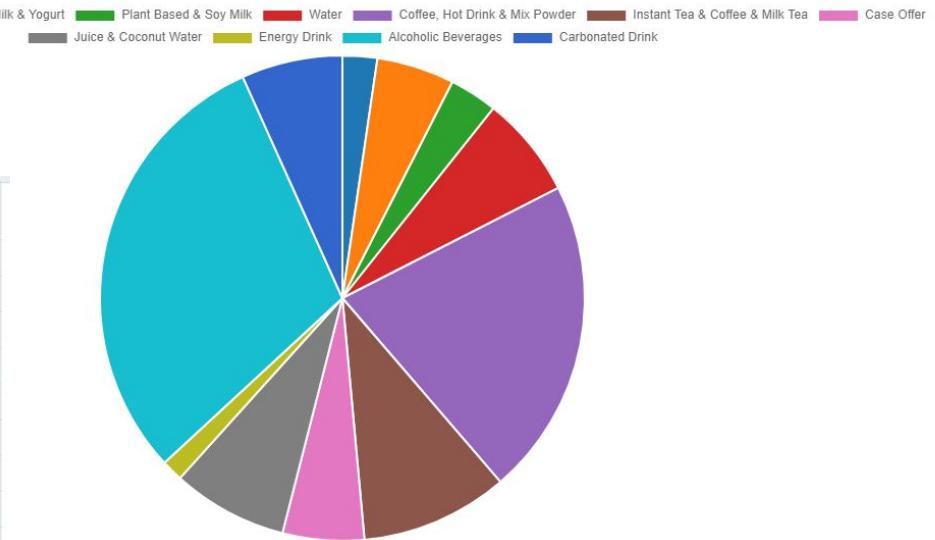
Graph Type: Pie Chart

Label: category

Value: no_of_product

| | category | no_of_product |
|----|---------------------------------|---------------|
| 1 | Herbal & Healthy Drink | 55 |
| 2 | Milk & Yogurt | 123 |
| 3 | Plant Based & Soy Milk | 76 |
| 4 | Water | 161 |
| 5 | Coffee, Hot Drink & Mix Powder | 503 |
| 6 | Instant Tea & Coffee & Milk Tea | 234 |
| 7 | Case Offer | 129 |
| 8 | Juice & Coconut Water | 183 |
| 9 | Energy Drink | 34 |
| 10 | Alcoholic Beverages | 715 |
| 11 | Carbonated Drink | 160 |

Using pgAdmin 4 Graph Visualiser



Analyzing data

average price of sugar-free drinks in last ten days

Query Query History

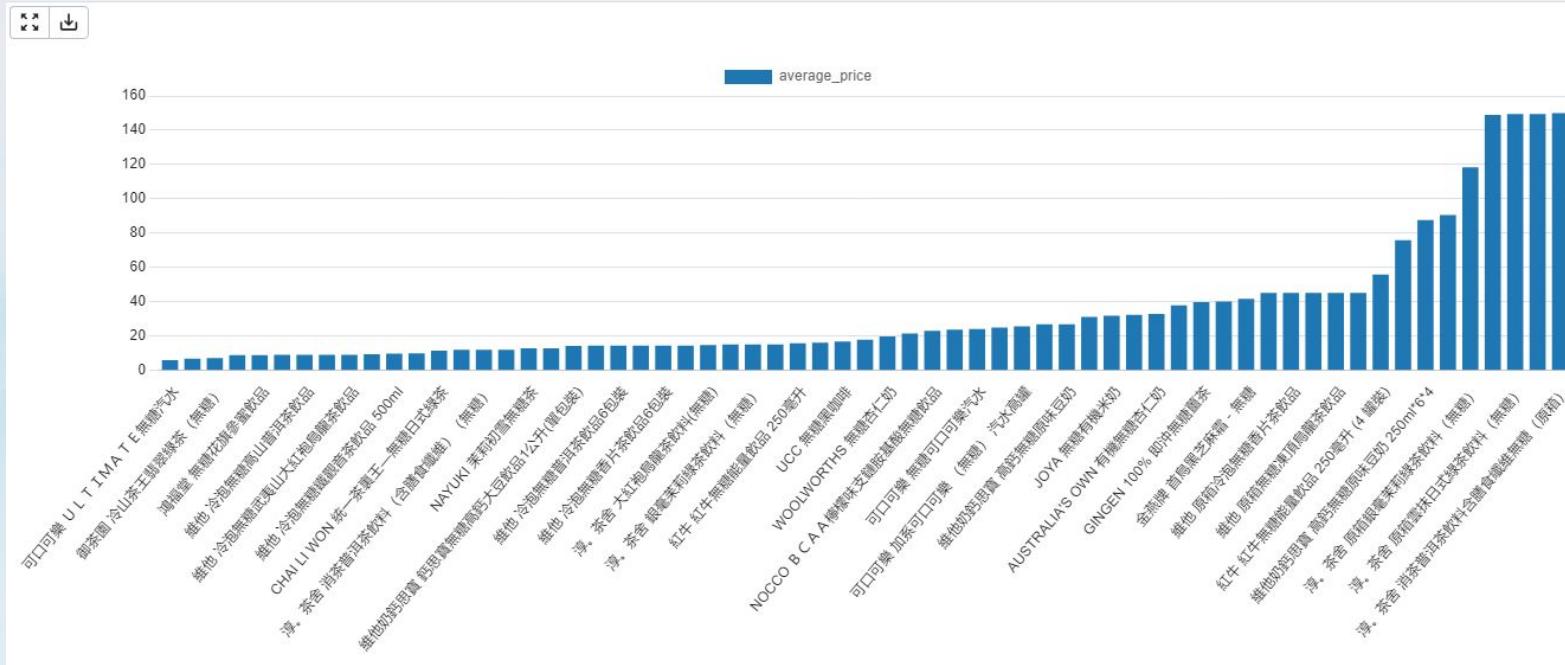
```
1 select name, avg(selling_price) as average_price from price
2 join product
3 using (link)
4 where name like '%無糖%'
5 and retailer = 'parknshop'
6 group by name
7 order by average_price;
```

| | name character varying (255) | average_price double precision |
|----|---------------------------------|-----------------------------------|
| 1 | 可口可樂 U L T I M A T E 無糖汽水 | 6 |
| 2 | ITAMACHI 麥茶 (無糖) | 6.90000000000001 |
| 3 | 御茶園 冷山茶王翡翠綠茶 (無糖) | 7.31666666666667 |
| 4 | 鴻福堂 雞漢果(無糖)飲品 | 8.90000000000002 |
| 5 | 鴻福堂 無糖花旗參蜜飲品 | 8.90000000000002 |
| 6 | 維他 冷泡無糖觀音茶飲品 500ml | 9.0818181818182 |
| 7 | 維他 冷泡無糖高山普洱茶飲品 | 9.15 |
| 8 | 維他 維他冷泡無糖香片茶飲品500ml | 9.15 |
| 9 | 維他 冷泡無糖武夷山大紅袍烏龍茶飲品 | 9.15 |
| 10 | 淳。茶舍 消夜茶飲料 (含膳食纖維) 無糖 | 9.44166666666668 |
| 11 | 維他 冷泡無糖鐵觀音茶飲品 500ml | 9.9 |
| 12 | 可口可樂 可口可樂加系可口可樂 (無糖) | 9.91666666666666 |
| 13 | CHAI LI WON 統一茶裏王-無糖日式綠茶 | 11.56666666666667 |
| 14 | 淳。茶舍 裏抹TM日式綠茶飲料(無糖) | 12.09583333333331 |
| 15 | 淳。茶舍 消夜普洱茶飲料 (含膳食纖維) (無糖) | 12.09583333333331 |
| 16 | 淳。茶舍 銀毫茉莉綠茶飲料(無糖) | 12.09583333333331 |
| 17 | NAYUKI 茉莉初雪無糖茶 | 12.9 |
| 18 | NAYUKI 邂逅紅茶無糖茶 | 12.9 |
| 19 | 維他奶奶思寶 鈣思寶無糖高鈣大豆飲品1公升(單包裝) | 14.25 |
| 20 | 維他 冷泡無糖凍頂烏龍茶飲品 | 14.5 |
| 21 | 維他 冷泡無糖普洱茶飲品6包裝 | 14.5 |
| 22 | 維他 冷泡無糖綠茶檸檬茶飲品 | 14.5 |
| 23 | 維他 冷泡無糖香片茶飲品6包裝 | 14.5 |
| 24 | 維他 冷泡無糖鐵觀音飲品 | 14.5 |
| 25 | 淳。茶舍 大紅袍烏龍茶飲料(無糖) | 14.75 |
| 26 | 淳。茶舍 淳。茶舍 大紅袍烏龍茶-無糖 | 15.16666666666666 |
| 27 | 淳。茶舍 銀毫茉莉綠茶飲料 (無糖) | 15.16666666666666 |
| 28 | 道地 無糖烏龍茶1.5升裝 | 15.25833333333335 |
| 29 | 紅牛 紅牛無糖能量飲品 250毫升 | 15.90000000000004 |
| 30 | 百事 可樂無糖 4 瓶裝 | 16.23333333333338 |
| 31 | UCC 無糖黑咖啡 | 16.90000000000002 |
| 32 | UCC 聰人無糖咖啡 | 18 |

Total rows: 63 of 63 Query complete 00:00:00.122

Analyzing data

average price of sugar-free drinks in last ten days



Challenge

Future Vision

- ❖ Enhanced Data Crawling
 - Refine the data crawler
 - Incorporate machine learning
 - Explore additional data sources
- ❖ Advanced Analytics
 - Develop advanced analytics
 - Implement predictive analytics
- ❖ User-Friendly Interface
 - Technical and less technical users
 - Drag-and-drop functionality, customizable dashboards
- ❖ Collaboration and Reporting:
 - Allow users share insights, reports, and findings
 - Automated report to PDF/PowerPoint/email.
- ❖ Integration with External Tools:
 - Tableau or Power BI more visualizations and storytelling.
 - Enable API share to other analytics tools.
- ❖ Continuous Feedback and Improvement:
 - Regularly collect feedback.
 - Conduct user testing sessions.

Thank You

Q&A Session