

## Importing libraries

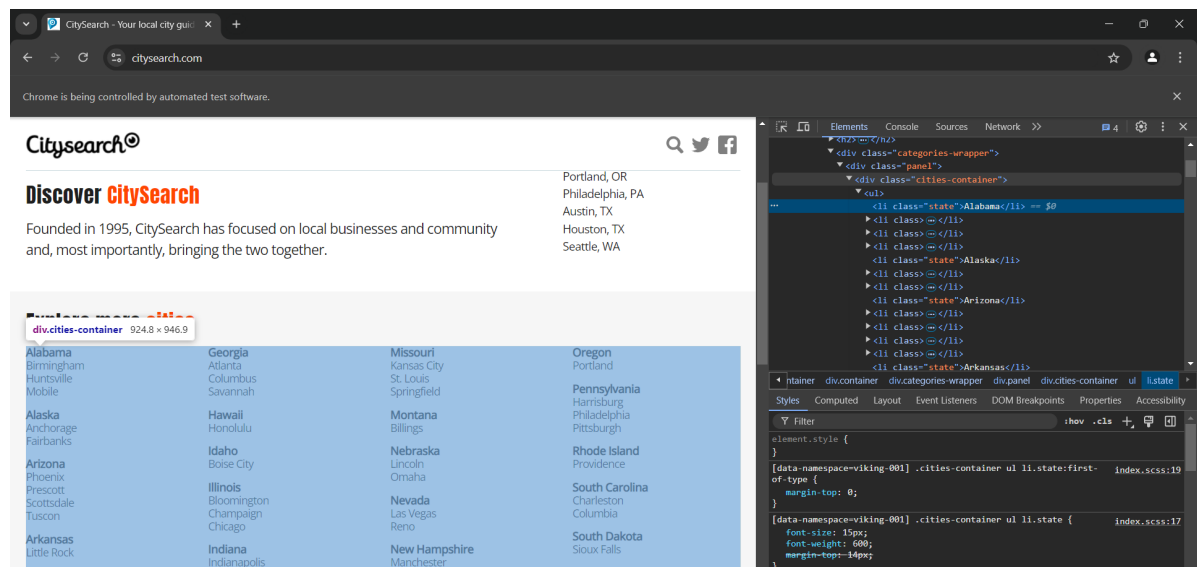
```
In [ ]: 1 import pandas as pd
        2 import time
        3 import re
        4
        5 from selenium import webdriver
        6 from selenium.webdriver.common.by import By
        7 from selenium.webdriver.support.ui import WebDriverWait
        8 from selenium.webdriver.support import expected_conditions as EC
        9 from selenium.common.exceptions import TimeoutException
```

## Navigating to the main page of CitySearch

```
In [ ]: 1 driver = webdriver.Chrome()
        2 driver.get("https://www.citysearch.com/")
```

## Extracting the links to individual cities

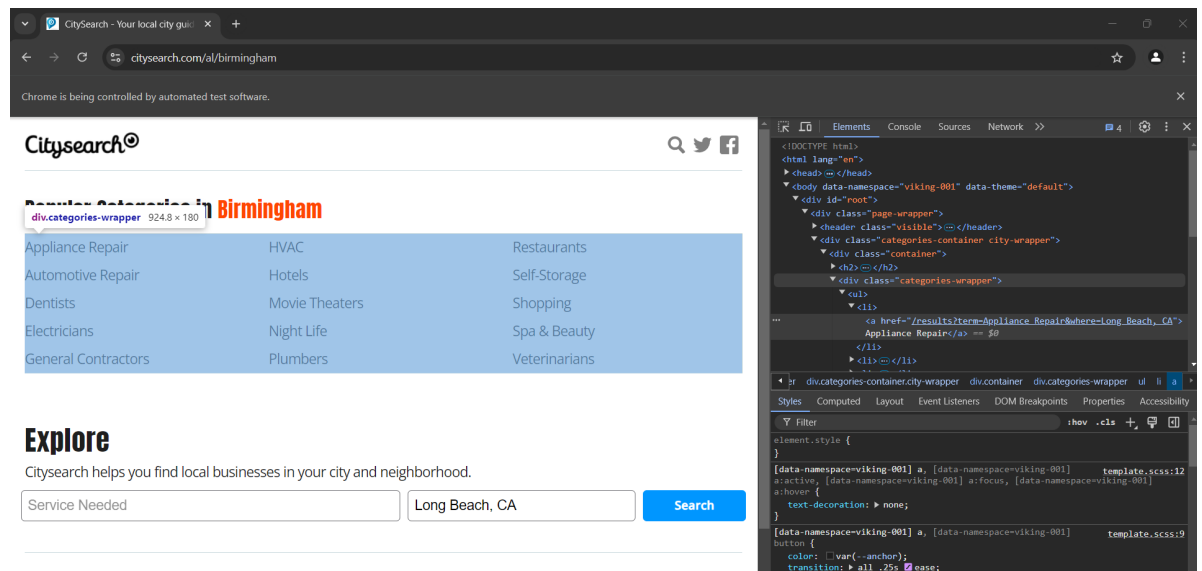
```
In [ ]: 1 container = driver.find_element(By.CSS_SELECTOR, "div.cities-container")
2 cities = container.find_elements(By.CSS_SELECTOR, "li:not([class*='sta')")
3
4 city_links = [city.get_attribute("href") for city in cities]
5 state, city = re.search("(?<=.com\\/).*", city_links[0]).group().split
```



## Navigating to a city link and gathering popular jobs

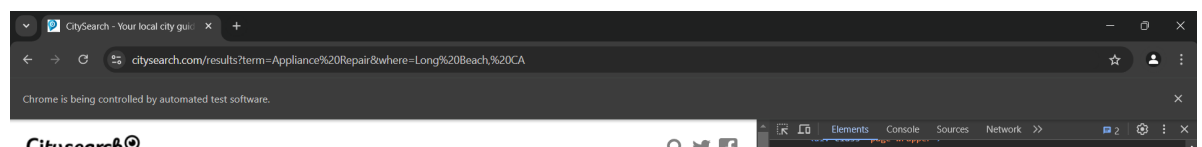
If we have keywords of specific industries we're interested in, I can iterate over them instead of iterating over popular industries. Also, if we have a list of states or cities we're interested in, I can also iterate over those.

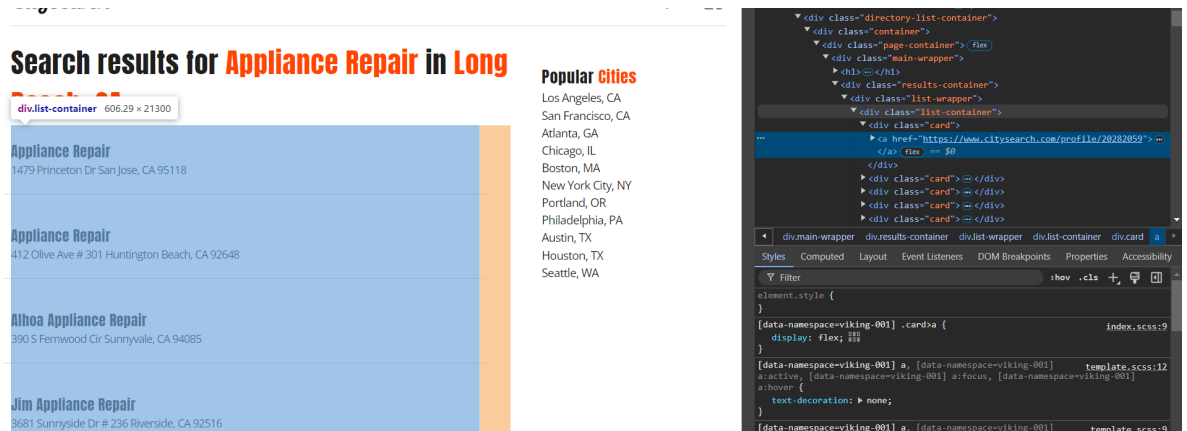
```
In [ ]: 1 driver.get(city_links[0]) # as an example will be going through the job
2 try:
3     elem = WebDriverWait(driver, 10).until(EC.presence_of_element_located((By.CSS_SELECTOR, 'div.categories-wrapper')))
4 except TimeoutException:
5     print("Timed out waiting for page to load")
6
7 popular_jobs = driver.find_elements(By.CSS_SELECTOR, 'div.categories-wrapper')
8 popular_jobs_links = [job.get_attribute("href") for job in popular_jobs]
```



## Navigating to first popular job and extracting links to jobs

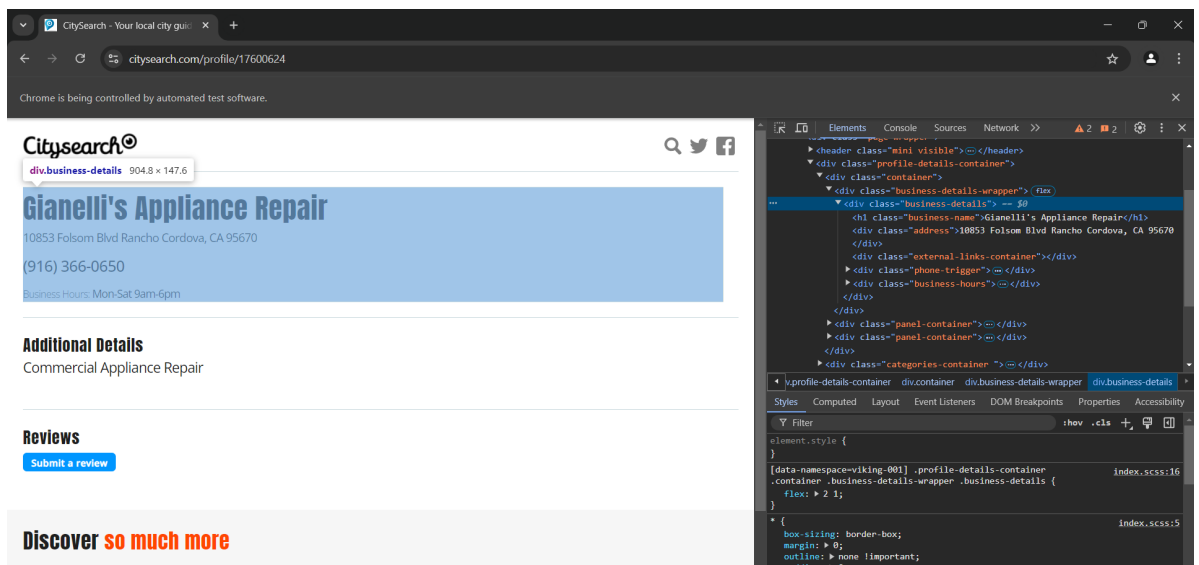
```
In [ ]: 1 driver.get(popular_jobs_links[0])
2
3 try:
4     elem = WebDriverWait(driver, 10).until(EC.presence_of_element_located((By.CSS_SELECTOR, 'div.list-container')))
5 except TimeoutException:
6     print("Timed out waiting for page to load")
7
8 job_cards = driver.find_elements(By.CSS_SELECTOR, 'div.list-container')
9 job_cards_links = [job.get_attribute("href") for job in job_cards]
```





## Scraping job description

```
In [ ]: 1 business_list = []
2
3 for i in range(0, 5): #5 should be job_cards_link's length when implem
4
5     driver.get(job_cards_links[i])
6
7     try:
8         elem = WebDriverWait(driver, 10).until(EC.presence_of_element_
9     except TimeoutException:
10         print("Timed out waiting for page to load")
11
12     # not sure if all business have all their contact info so creating
13     business_details = driver.find_elements(By.CSS_SELECTOR, 'div.busi
14
15     business_details_dict = {
16         entry.get_attribute("class"): entry.text
17         for entry in business_details
18     }
19
20     business_list.append(business_details_dict)
21     time.sleep(3)
```



## Converting to dataframe, renaming columns and exporting to csv

```
In [ ]: 1 df = pd.DataFrame.from_dict(business_list)
2
3 df.rename(columns={
4     "business-name": "business name",
5     "external-links-container": "external link",
6     "phone-trigger": "phone number",
7     "business-hours": "business hours"
8 }, inplace=True)
9
10 df.to_csv(f'./{state}_{city}.csv', index=False) # example al_birmingham
```

```
In [ ]: 1 driver.quit()
```

## Possible Improvements and Changes

Depending on the company's needs, I can store these values elsewhere instead of a CSV. Possibly in MongoDB or an SQL database.

When scraping large amounts of data, the current script may run into memory issues. During full implementation, I'll refactor the script into something more modular or OOP. Selenium has something called Page Object Model (POM). I'm not too familiar with POM but I am more than willing to try!