# CitySearch Web Scraping

## Implementation (8/20/24)

### Importing libraries

```python
import pandas as pd
import time
import re

from selenium import webdriver
from selenium.webdriver.common.by import By
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
from selenium.common.exceptions import TimeoutException
```
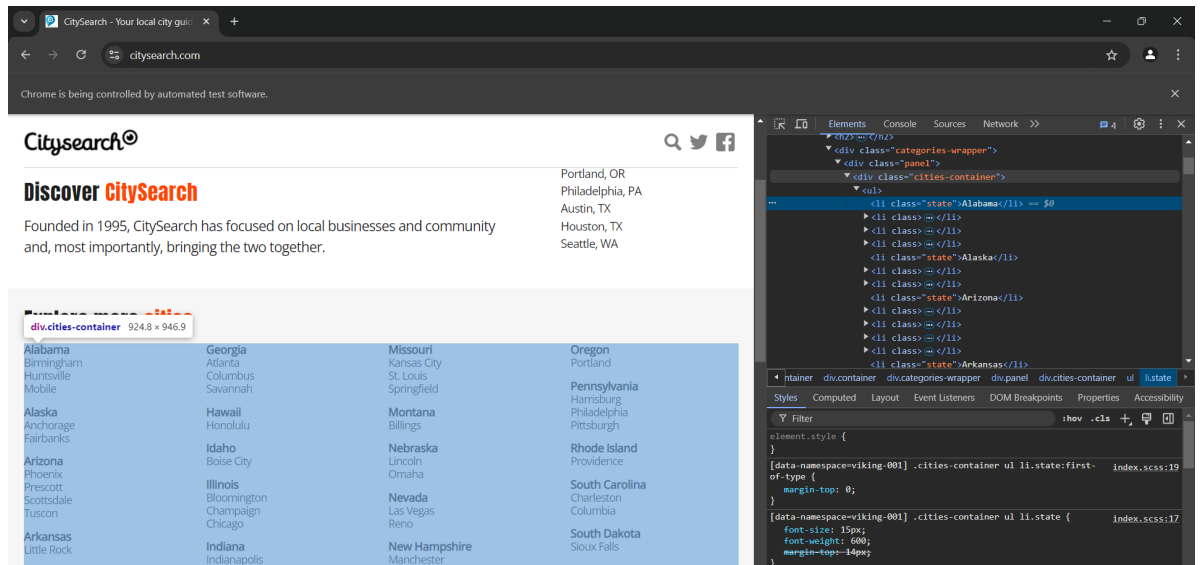
### Navigating to the main page of CitySearch

```python
driver = webdriver.Chrome()
driver.get("https://www.citysearch.com/")
```

### Extracting the links to individual cities

```python
container = driver.find_element(By.CSS_SELECTOR, "div.cities-container
cities = container.find_elements(By.CSS_SELECTOR, "li:not([class*='sta

city_links = [city.get_attribute("href") for city in cities]
state, city = re.search("(?<=\.com\/).*", city_links[0]).group().split
```
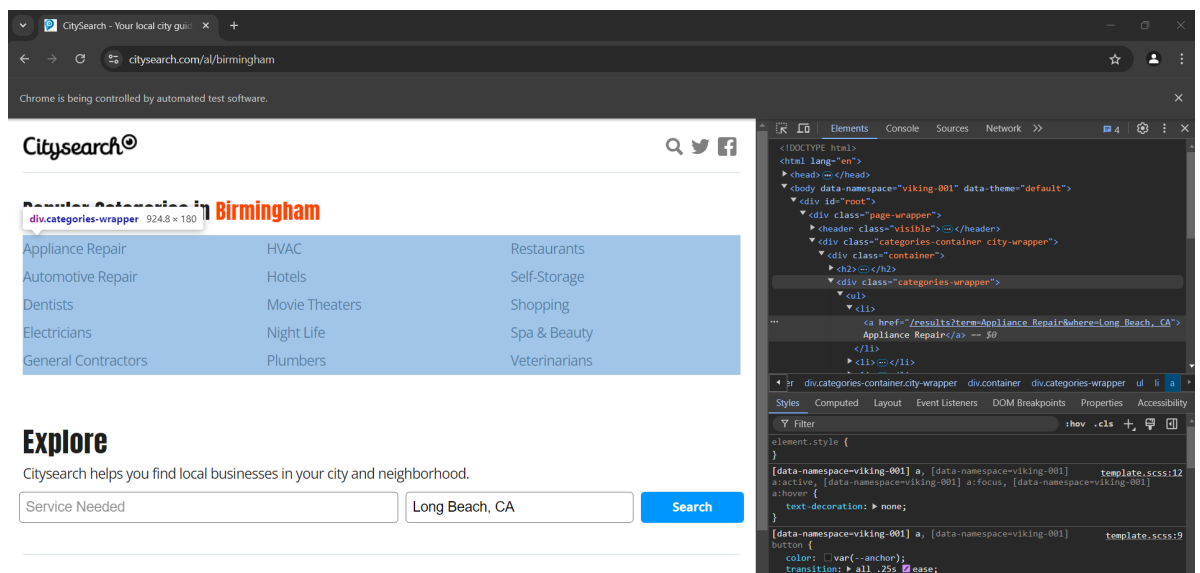
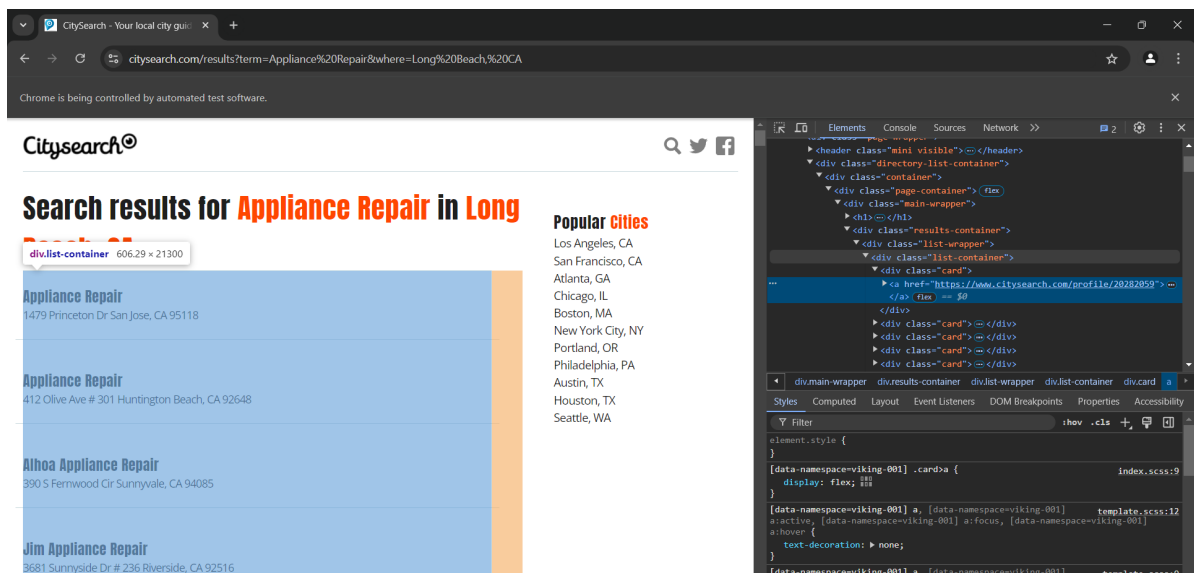## Navigating to a city link and gathering popular jobs

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 jo
2  try:
3      elem = WebDriverWait(driver, 10).until(EC.presence_of_element_loca
4  except TimeoutException:
5      print("Timed out waiting for page to load")
6
7  popular_jobs = driver.find_elements(By.CSS_SELECTOR, 'div.categories-w
8  popular_jobs_links = [job.get_attribute("href") for job in popular_job
```

## 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_loca
        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]
```
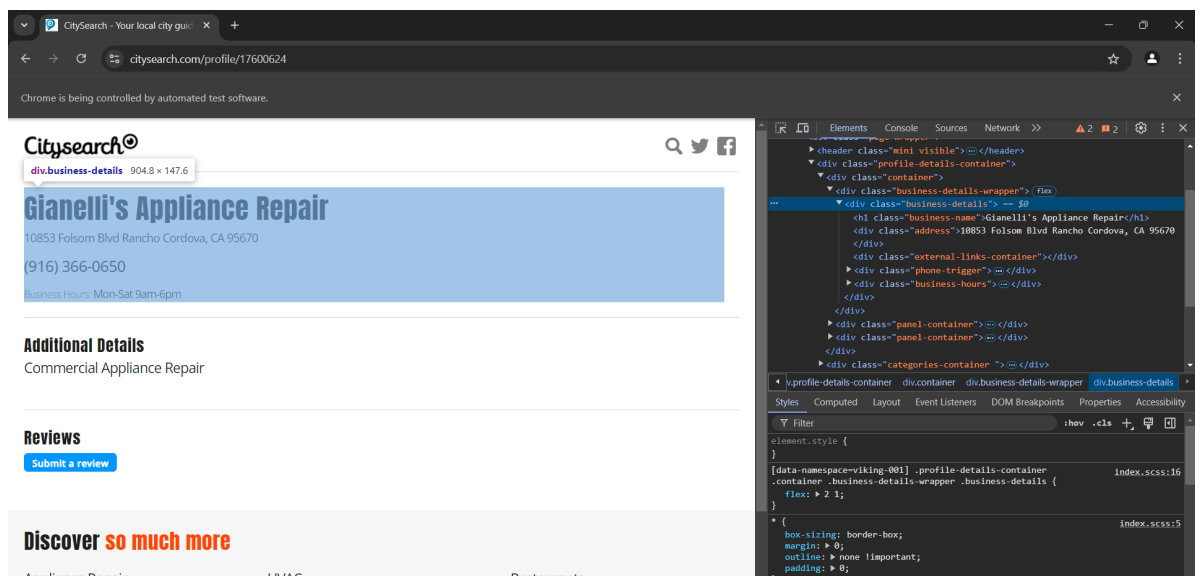


## Scraping job description

```python
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_birmingha
```

```
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!

# Implementation (8/26/24)

## Importing Libraries

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

## Key for converting state name to acronym

```python
us_state_to_abbrev = {
    "Alabama": "AL",
    "Alaska": "AK",
    "Arizona": "AZ",
    "Arkansas": "AR",
    "California": "CA",
    "Colorado": "CO",
    "Connecticut": "CT",
    "Delaware": "DE",
    "Florida": "FL",
    "Georgia": "GA",
    "Hawaii": "HI",
    "Idaho": "ID",
    "Illinois": "IL",
    "Indiana": "IN",
    "Iowa": "IA",
    "Kansas": "KS",
    "Kentucky": "KY",
    "Louisiana": "LA",
    "Maine": "ME",
    "Maryland": "MD",
    "Massachusetts": "MA",
    "Michigan": "MI",
    "Minnesota": "MN",
    "Mississippi": "MS",
    "Missouri": "MO",
    "Montana": "MT",
    "Nebraska": "NE",
    "Nevada": "NV",
    "New Hampshire": "NH",
    "New Jersey": "NJ",
    "New Mexico": "NM",
    "New York": "NY",
    "North Carolina": "NC",
    "North Dakota": "ND",
    "Ohio": "OH",
    "Oklahoma": "OK",
    "Oregon": "OR",
    "Pennsylvania": "PA",
    "Rhode Island": "RI",
    "South Carolina": "SC",
    "South Dakota": "SD",
    "Tennessee": "TN",
    "Texas": "TX",
    "Utah": "UT",
    "Vermont": "VT",
    "Virginia": "VA",
    "Washington": "WA",
    "West Virginia": "WV",
    "Wisconsin": "WI",
    "Wyoming": "WY",
    "District of Columbia": "DC",
    "American Samoa": "AS",
    "Guam": "GU",
    "Northern Mariana Islands": "MP",
```

```
56        "Puerto Rico": "PR",
57        "United States Minor Outlying Islands": "UM",
58        "U.S. Virgin Islands": "VI",
59  }
60
```

## Importing .xlxs file and filling in merged cells with previous value

In [ ]:  ▶|
```
1  df = pd.read_excel("google_maps_keywords.xlsx")
2  df.loc[:, ["Country", "State"]] = df.loc[:, ["Country", "State"]].ffil
```

## Opening CitySearch

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

## Finding container for links of all cities and saving it to a variable

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  city_links = [city.get_attribute("href") for city in cities]
4  state, city = re.search("(?<=\.com\/).*", city_links[0]).group().split
5
```

## Grouping dataframe by country and state

During implementation, I would look over all the countries and state. This would be outer most loop. Until implementation, I have hard coded the script to only loop over the cities in the US.

In [ ]:  ▶|
```
1  grouped_df = df.groupby("Country")
2  grouped_countries = grouped_df.get_group("United States") #todo during
3  grouped_states = grouped_countries.groupby("State")
4
5  states = grouped_states.groups.keys()
```

## Looping for the states, cities, and industries in that order.

Some cities didn't have any job postings for some industries, so I've added a TryExcept block to cover for those situations. Currently, it's only looping for the first 5 job postings, but during implementation, the loop would continue until it's extracted everything.

```python
for state in states:
    state_cities = grouped_states.get_group(state)["City"]
    state_industries = grouped_states.get_group(state)["Industry"]
    state_abbrev = us_state_to_abbrev[state]

    for city in state_cities:
        business_list = []

        if (pd.isnull(city)): continue

        for industry in state_industries:
            if (pd.isnull(industry)): continue

            url = f"https://www.citysearch.com/results?term={industry.
            print("-------------------------", url, "------------------
            driver.get(url)

            # extracting all links to jobs in this category
            try:
                elem = WebDriverWait(driver, 10).until(
                    EC.presence_of_element_located((By.CSS_SELECTOR, "
            except TimeoutException:
                print("Timed out waiting for page to load: Most likely

            job_cards = driver.find_elements(By.CSS_SELECTOR, "div.lis
            job_cards_links = [job.get_attribute("href") for job in jo

            # visiting each job link for the current industry and scra

            if (len(job_cards_links) == 0): continue

            # for i in range(len(job_cards_links)): # todo uncomment t
            for i in range(5):
                print(job_cards_links[i])
                driver.get(job_cards_links[i])

                try:
                    elem = WebDriverWait(driver, 10).until(
                        EC.presence_of_element_located((By.CSS_SELECTO
                except TimeoutException:
                    print("Timed out waiting for page to load")

                # not sure if all business have all their contact info
                business_details = driver.find_elements(By.CSS_SELECTO

                business_details_dict = {
                    entry.get_attribute("class"): entry.text
                    for entry in business_details
                }
                business_details_dict["industry"] = industry

                business_list.append(business_details_dict)
                time.sleep(3)

        df = pd.DataFrame.from_dict(business_list)
```

```
56
57          df.rename(columns={
58              "business-name": "business name",
59              "external-links-container": "external link",
60              "phone-trigger": "phone number",
61              "business-hours": "business hours"
62          }, inplace=True)
63
64          df.to_csv(f'./{state_abbrev.lower()}_{city.lower()}.csv', inde
65
```

In [ ]: ▶  1  driver.quit()

# Results



**I think it's ready for full implementation!**