# Data and Artificial Intelligence
# Cyber Shujaa Program

## Week 1 Assignment
## Web Scraping and Data Handling in Python

**Student Name:** John Brown Ouma

**Student ID:** CS-DA01-25030.

## Introduction

This report documents the completion of a web scraping assignment that involved extracting structured data from the "Forms and Search" page on scrapethissite.com. The goal was to gather hockey team statistics, organize them into a structured format using pandas, and export the data to a CSV file for further analysis. The entire scraping process was implemented using Python on Google Colab, utilizing key libraries including requests for making HTTP requests, BeautifulSoup for parsing HTML, and pandas for data management and export.

The objectives of the assignment were:
1. Practical Python coding on Jupiter Notebooks hosted on Google Colab
2. Use requests and BeautifulSoup to extract data from a web page.
3. Parse and clean the extracted data.
4. Store structured data into a Pandas DataFrame.
5. Export the final dataset to a .csv file.

## Tasks Completed

### 1. Setting Up the Environment

Before web scraping, it's essential to import the libraries required for the task and establish a connection to the target webpage. In this step:

- **requests** fetches the webpage's HTML content.

- **BeautifulSoup** is responsible for parsing the HTML so it can be navigated and queried.

- **pandas** will later help structure the scraped data into a tabular format.

*Figure 1: Importing libraries and setting up the target URL*

This ensures the environment has all tools needed for making HTTP requests, parsing HTML, and handling data.

## 2. HTML Parsing and Table Extraction

Once the page's content is fetched, we parse it into a **BeautifulSoup object**, allowing us to interact with the page's structure and locate elements of interest — in this case, the **hockey statistics table**.



*Figure 2: Parsing HTML*



*Figure 2.1 finding the target table*

This transforms the raw HTML string into a navigable object and finds the specific table containing the data by searching for the table element with class table.

## 3. Extracting Column Headers

Before collecting data, it's crucial to identify the table's **column titles** to accurately structure our data. This is done by locating all <th> tags within the table.



*Figure 3: Extracting and displaying table column headers*

This retrieves clean, human-readable column names which will act as the DataFrame's headers, ensuring our dataset is well-organized.

## 4. Creating the DataFrame and Populating with Data

With headers identified, the next step is to **initialize an empty DataFrame** using those headers and then **loop through each row in the table body** (skipping the header row) to extract data.



*Figure 4: Populating the DataFrame with scraped data*

This systematically extracts each cell's content, cleans it (removes leading/trailing whitespace), and appends it to the DataFrame — one row at a time.

## 5. Final DataFrame Inspection

After populating the DataFrame, it's important to **review the resulting table** to ensure data integrity and verify the scraping process was successful.



*Figure 5:DataFrame Inspection*

Quickly displays the entire structured dataset within the Colab notebook, providing an opportunity to visually inspect for missing values or irregularities.

## 6. Exporting to CSV

Finally, we export the clean, structured dataset to a **CSV file** for future use or sharing.
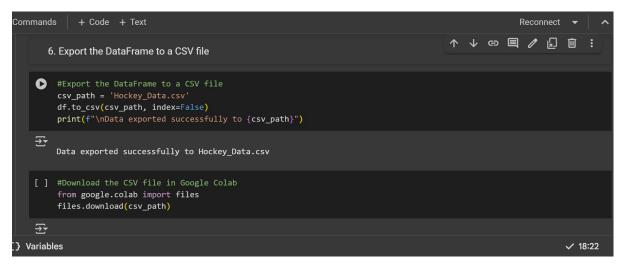


*Figure 6: Exporting to CSV and downloading the CSV*

CSV is a versatile, widely supported file format suitable for data analysis, visualization, or integration into other tools or projects.

**Link to Code:**

https://colab.research.google.com/drive/1agOHXQMxoHCiqEyRxSHWb9IvQHRjpUWJ?usp=sharing

## Conclusion

This assignment successfully demonstrated the process of web scraping using Python. I was able to:

1. Retrieve web page content using the Requests library

2. Parse and extract specific data using BeautifulSoup

3. Structure and clean the data using pandas

4. Export the final dataset to a CSV file

The complete Google Colab notebook can be accessed here: https://colab.research.google.com/drive/1agOHXQMxoHCiqEyRxSHWb9IvQHRjpUWJ?usp=sharing

Through this exercise, I gained practical experience in web scraping techniques and data organization, which are valuable skills for data collection and analysis tasks.