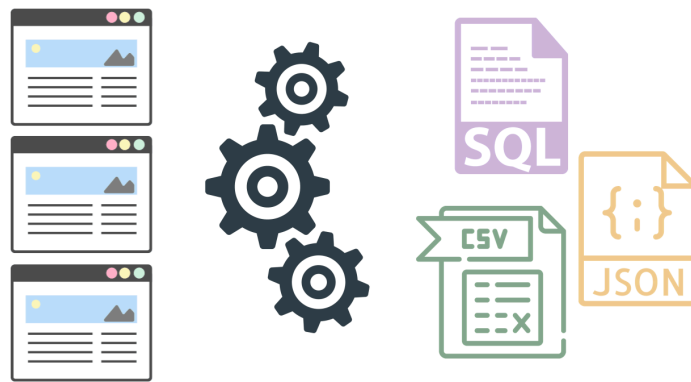


# Data Scraping and Cleaning Methodologies for the StudyBuddy Hub Project



## Work Distribution:

Scraped Courses, Colleges, Resources Data

Sarah Kayembe

Scraped Quizzes Data

Vivek Reddy Bhimavarapu

Cleaned Courses, Colleges, Quizzes Data

Rise Akizaki

Cleaned Resources Data

Jacob Craig

# Table of Contents

## [Table of Contents](#)

### [Data Scraping](#)

#### [College Table](#)

- [1. Objective and Justification](#)
- [2. Methodology and Tools](#)
- [3. Data Source and Scraping Strategy](#)
  - [1. Data Source](#)
  - [2. Scraping Strategy](#)
- [4. Error Handling](#)
- [5. Output](#)

#### [Course table](#)

- [1. Objective and Justification](#)
- [2. Methodology and Tools](#)
- [3. Data Source and Scraping Strategy](#)
  - [1. Data Source](#)
  - [2. Scraping Strategy](#)
- [4. Error Handling](#)
- [5. Output](#)

#### [Resources Table](#)

- [1. Objective and Justification](#)
- [2. Methodology and Tools](#)
- [3. Data Source and Scraping Strategy](#)
  - [1. GeeksforGeeks – Data Structures Section](#)
  - [2. Khan Academy – Biology Units](#)
  - [3. W3Schools \(Python & Java Tutorials\)](#)
  - [4. Criminal Law \(Multi-Source Scrape\)](#)
- [4. Error Handling and Validation Techniques](#)
  - [1. Network Errors:](#)
  - [2. Data Validation:](#)
  - [3. File Validation:](#)
- [5. Output](#)

### [Data Cleaning](#)

#### [usm\\_courses](#)

- [1. Algorithm For Cleaning](#)
- [2. Saving Results](#)

#### [us\\_colleges](#)

- [1. Algorithm For Cleaning](#)
- [2. Saving Results](#)

#### [quiz\\_data](#)

1. Questions Data Frame

Algorithm For Cleaning / Saving Results

2. Answers Data Frame

Algorithm For Cleaning

3. Saving Results

course\_resources

1. Objectives and Justification

2. Methodology & Tools

3. Cleaning Algorithm

4. Cleaned DataFrame

5. Validation and Output

Appendix

This report explains how we pulled and cleaned the data to make reliable academic datasets for the StudyBuddy Hub project. Our main aim was to automatically gather good information on colleges, courses, and resources so that way, we skip the mistakes that happen when people type data in by hand.

We used Python and tools like requests, BeautifulSoup, and pandas to grab the data from solid online places, including Wikipedia and college catalogs. After collecting it, we cleaned, checked, and removed repeats to make sure everything was accurate, looked consistent, and would work easily with the project's database.

# Data Scraping

## College Table

### 1. Objective and Justification

The U.S. College Name Scraper was created to create a standardized database of American colleges and universities by extracting institution names from Wikipedia's "Lists of American universities and colleges."

The primary goal of this scraping task was to prevent inconsistent or duplicate college entries within the application. Without a verified list, users could manually input college names with spelling variations or abbreviations (e.g., "MIT," "Massachusetts Inst. of Tech," "Massachusetts Institute of Technology"), which would fragment data and incorrectly categorize students under different institutions.

By automating the collection of verified college names, the system ensures uniformity and accurate filtering.

### 2. Methodology and Tools

The data collection was executed using a Python script (documented in `scraping_college_names.ipynb`) utilizing the following libraries:

Tool	Purpose
Python requests	Handled HTTP requests to retrieve raw HTML content.

**Python BeautifulSoup**                      Parsed the HTML to locate and extract specific data elements.

**Python pandas**                      Structured the data, removed duplicates, sorted the final output, and managed CSV persistence.

### **3. Data Source and Scraping Strategy**

#### **1. Data Source**

The source was List of American Universities and Colleges on Wikipedia with the URL [https://en.wikipedia.org/wiki/Lists\\_of\\_American\\_universities\\_and\\_colleges](https://en.wikipedia.org/wiki/Lists_of_American_universities_and_colleges). The universities were listed within a U.S. state or territory in a tabular form under standardized classes.

#### **2. Scraping Strategy**

##### **Hub Page Crawling:**

- The scraper began by fetching the main hub page, which contains hyperlinks to all 50 states and territory-specific lists.
- Using BeautifulSoup, it located all <a> tags whose href attribute matched the pattern /wiki/List\_of\_colleges\_and\_universities\_in\_....
- Each link's absolute URL was reconstructed using urljoin, and the state name was extracted as the key identifier.

##### **Iterative State-Level Scraping:**

- For each state link, the scraper requested the page and searched for all <table class="wikitable"> elements.
- The table with the most rows was prioritized, as it generally contained the main list of institutions.
- The HTML for the selected table was then parsed into a pandas DataFrame using pd.read\_html() for structured extraction.

##### **Column Identification:**

- Because Wikipedia's tables vary slightly between states, the scraper dynamically identified the correct column by testing for one of the following headers: ['School', 'Institution', 'University', 'College', 'Name', 0].
- The first valid match was used to extract college names.

#### Data Cleaning:

- Bracketed citation markers (e.g., [1], [a]) were removed from all names using string operations.
- Whitespace and special characters were stripped.
- Each entry was stored with its corresponding state for context.

#### Politeness and Retry Policy:

- All HTTP requests were made with a realistic browser User-Agent header to avoid being flagged as a bot.
- Each page was fetched up to three times in case of network or timeout errors.
- Randomized polite delays (`time.sleep(1 + random.random())`) were applied to prevent IP throttling.

### 4. Error Handling

- **Retry Logic:** Each HTTP request included up to three retries with incremental delays to handle temporary connectivity issues.
- **Graceful Failures:** States that could not be fetched or parsed after three attempts were skipped with a warning message, allowing the script to continue without interruption.
- **Fallback Safety:** If no valid table was found on a page, the scraper safely returned an empty list for that state rather than terminating execution.

### 5. Output

- **File Output:** `us_colleges.csv`, `us_colleges_clean.csv`

Column	Description
State	The U.S. state where the institution is located

**College Name**      Official name of the college or university as listed on Wikipedia

Entries were cleaned and deduplicated based on College Name + State, ensuring only one standardized instance of each institution exists.

## Course table

### 1. Objective and Justification

The goal was to build a standardized database of higher education course identifiers and descriptions. This was done by scraping the course catalog of a single institution, the University of Southern Maine (USM), to establish a controlled, limited scope of authentic, real-world course data for system testing and demonstration purposes.

The scraper uses Python libraries such as requests, BeautifulSoup, and pandas to navigate, parse, and structure the course data. The output is stored in a CSV file (usm\_courses.csv) within a data directory for further integration.

### 2. Methodology and Tools

The data collection was executed using a Python script (documented in scraping\_courses.ipynb) utilizing the following libraries:

Tool	Purpose
Python requests	Fetches the course catalog's pagination pages.
Python BeautifulSoup	Parsed the HTML of each catalog page to locate and extract course tables.
Python pandas	Consolidated the data, removed duplicates, and managed the final CSV output.

### 3. Data Source and Scraping Strategy

#### 1. Data Source

The source was the University of Southern Maine's (USM) public online course catalog with the base URL as <https://catalog.usm.maine.edu/content.php>. The content type was dynamic HTML pages displaying course listings filtered through query parameters (e.g, filter[cpage] for pagination).

## **2. Scraping Strategy**

### **Pagination Detection:**

- To determine the total number of catalog pages, the scraper first probed a known content page (Page 2).
- Using a regular expression search on strings containing the pattern "Page X of Y", it dynamically extracted the total page count (Y).
- If pagination text could not be located, the program defaulted to 18 pages as a safe fallback.

### **Iterative Extraction:**

- For each page, the scraper constructed HTTP GET requests with the appropriate query parameters.
- It then parsed the resulting HTML to find all <a> tags whose href contained the substring "preview\_course\_nopop.php", indicating valid course entries.
- The visible text of each link (formatted as "CODE - COURSE NAME") was split into course\_code and course\_name.
- Each record was appended to a list with the associated course\_link and page number for traceability.

### **Human Validation via Inspection:**

- Prior to coding, multiple pages were manually inspected using Chrome DevTools to understand structural patterns and ensure the consistency of link and text placement.
- It was confirmed that the structure of <a> tags and link patterns remained consistent across all catalog pages.

### **Politeness Policy:**



- A 0.5-second delay (`time.sleep(0.5)`) between requests was implemented to minimize load on the USM server and comply with polite scraping norms.

## 4. Error Handling

- **Request Exceptions:** Every HTTP call was wrapped in a try-except block to catch network errors and unexpected responses (`RequestException`).
- **Pagination Failures:** If a page failed to load or no courses were found mid-sequence, the scraper automatically stopped pagination to avoid redundant requests.
- **Fallback Safety:** Default values were used for both total page count and missing course metadata to ensure uninterrupted execution and dataset completeness.

## 5. Output

- **File Output:** `usm_courses.csv`

Column	Description
<code>course_code</code>	Official course (e.g., COS 170)
<code>course_name</code>	Full course title as listed in the catalog (e.g., Structured Programming Lab)
<code>course_link</code>	URL to the individual course details page
<code>page_number</code>	Catalog page number from which the course was extracted

Duplicate entries were removed based on the combined fields `course_code` + `course_name` to ensure data uniqueness.

# Resources Table

## 1. Objective and Justification

The purpose of this scraping module was to populate the studyBuddy Hub resource table with verified, structured learning materials relevant to specific courses in the form of links, articles, pdfs, videos, unit overviews, etc.

**Crucial Strategy Shift:** The initial approach using broad web crawling was discontinued due to the inconsistent quality and limited relevance of the retrieved content. Instead, a more focused strategy was adopted—targeting specific, pre-vetted educational websites to ensure that all collected resources are accurate and directly map to course content.

## 2. Methodology and Tools

The data collection was executed using a Python script (documented in `scraping_resources.ipynb`) utilizing the following libraries:

Tool	Purpose
Python requests	Sent HTTP requests and retrieved HTML page content from the target websites.
BeautifulSoup (bs4)	Parsed HTML and extracted specific elements like titles ( <code>&lt;h1&gt;/&lt;b&gt;</code> ), descriptions ( <code>&lt;p&gt;</code> ), and URLs ( <code>href</code> ).
pandas	Structured the extracted data into tabular form, handled duplicates, standardized column names, and appended records to the master CSV file ( <code>course_resources_scraped.csv</code> ).

## 3. Data Source and Scraping Strategy

The scraping process was designed to collect structured, course-specific educational resources from diverse and reputable open-access learning sites. Each source required manual inspection of its HTML structure to identify consistent tag hierarchies and reliable selectors for extracting topics, descriptions, and URLs. All scripts were implemented in Python using the requests, BeautifulSoup, and pandas libraries, with outputs consolidated into a unified CSV file (`course_resources_scraped.csv`) for integration into the resource table.

### 1. GeeksforGeeks – Data Structures Section

**Purpose:** Extract structured computer science topics with corresponding short descriptions and resource URLs.

**URLs:**

<https://www.geeksforgeeks.org/dsa/introduction-to-data-structures/>

<https://www.geeksforgeeks.org/dsa/introduction-to-algorithms/>

<https://www.geeksforgeeks.org/web-tech/web-technology/>

<https://www.montclair.edu/computer-science-education/wp-content/uploads/sites/253/2024/02/3-5-8.1.5-Algorithms-Programming.pdf>

**Inspection Process:**

- Used Chrome DevTools to inspect the structure of each page and verify the location of core tags such as <a rel="noopener">, <h2>, and <p dir="ltr">.
- Identified consistent placement of bolded topic titles inside <a> elements and accompanying text in adjacent <p> tags.

### Scraping Logic:

- Parsed HTML using BeautifulSoup and iterated through the detected topic links.
- Extracted the topic name, description, and hyperlink.
- Applied error handling to skip missing or malformed entries.

### Data Labeling:

- Each record was stored with ResourceType = 'LINK' and tagged to its corresponding CourseName (e.g., *Data Structures, Algorithms in Programming, Web Applications Development*).
- Additional PDF entries (e.g., *Algorithms Programming Concepts*) were manually appended to ensure coverage of offline resources.

## 2. Khan Academy – Biology Units

**Purpose:** Retrieve unit titles and summaries for biological sciences courses.

**URLs:** <https://www.khanacademy.org/science/ap-biology/chemistry-of-life>

<https://www.khanacademy.org/science/ap-biology/cell-structure-and-function>

<https://www.khanacademy.org/science/ap-biology/cellular-energetics>

<https://www.khanacademy.org/science/ap-biology/cell-communication-and-cell-cycle>

<https://www.khanacademy.org/science/ap-biology/heredity>

<https://www.khanacademy.org/science/ap-biology/gene-expression-and-regulation>

<https://www.khanacademy.org/science/ap-biology/natural-selection>

<https://www.khanacademy.org/science/ap-biology/ecology-ap>

<https://www.khanacademy.org/science/ap-biology/x16acb03e699817e9:simulations>

<https://www.khanacademy.org/science/ap-biology/worked-examples-ap-biology>

<https://www.khanacademy.org/science/ap-biology/ap-biology-standards-mappings>

#### **Inspection Process:**

- Reviewed the DOM hierarchy to confirm that unit identifiers were embedded in the URL slug rather than as structured headings.
- Verified that summary text appeared within <div> containers matching patterns like main-content or unit-description.

#### **Scraping Logic:**

- Extracted TopicName directly from URL segments using regular expressions to transform slugs into readable titles (e.g., /cellular-energetics → Cellular Energetics).
- Could not retrieve the first meaningful sentence of the page.

#### **Validation:**

- Used regular expressions to clean unwanted characters from titles.
- Applied checks for duplicates and malformed text.

### **3. W3Schools (Python & Java Tutorials)**

**Purpose:** Build comprehensive topic lists for programming languages.

**URLs:** <https://www.w3schools.com/java/>

<https://www.w3schools.com/python/>

#### **Inspection Process:**

- Inspected the sidebar navigation (id="leftmenuinnerinner") to locate the full set of tutorial links.

- Confirmed consistent structural patterns across both Python and Java tutorials.

#### Scraping Logic:

- Collected all .asp links from the navigation menu using urljoin to form absolute URLs.
- For each page, extracted the main <h1> title and the first paragraph as a description snippet.
- Used regular expressions to isolate the first full sentence for clarity.

#### Data Validation:

- Ensured every entry contained a valid title, snippet, and resource URL.
- Added a short delay (time.sleep()) between requests as a politeness measure.

### 4. Criminal Law (Multi-Source Scrape)

**Purpose:** Gather legal resource summaries from government, academic, and MOOC sites.

#### Sources:

- Maine Legislature:  
<https://legislature.maine.gov/statutes/17-a/title17-ach0sec0.html>
- American Public University System (APUS):  
<https://www.apu.apus.edu/area-of-study/security-and-global-studies/resources/what-is-criminal-law-and-why-does-it-matter/>
- edX: <https://www.edx.org/learn/criminal-law>

#### Inspection Process:

- Inspected each site individually due to structural variation.
- Located heading tags (<h1>, <h2>) and primary content containers (<article>, <main>, .entry-content, .main-content).

#### Scraping Logic:

- Extracted titles from visible headers or <title> tags as a fallback.
- Captured the first descriptive paragraph using regex to isolate the initial complete sentence.
- For edX, incorporated meta descriptions when available for cleaner summaries.

## 4. Error Handling and Validation Techniques

### 1. Network Errors:

All requests were wrapped in a try-except block to catch timeout, connection, or 404 errors.

**Example:**

```
try:
    response = requests.get(url, headers=HEADERS, timeout=10)
    response.raise_for_status()
except requests.exceptions.RequestException as e:
    print(f" * ERROR: Failed to fetch {url}. Reason: {e}")
```

### 2. Data Validation:

- Ignored links with missing text or invalid URLs.
- Filtered duplicates using Pandas .drop\_duplicates().
- Checked snippet lengths and trimmed overly long descriptions to 250-500 characters.

### 3. File Validation:

- Ensured all output columns matched the schema: ['CourseName', 'TopicName', 'ResourceURL', 'ContentSnippet', 'ResourceType'].

**Validated CSV creation only if data existed:**

```
if not df.empty:
    df.to_csv("course_resources_scraped.csv", index=False)
```

## 5. Output

Each individual script appended its results to a shared CSV file (course\_resources\_scraped.csv), creating a cumulative dataset of cross-disciplinary educational resources.

Column	Description
CourseName	Name of the associated course (e.g., “Data Structures”)
TopicName	Topic title extracted from the source page
ResourceURL	Direct link to the resource
ContentSnippet	Short description of the topic
ResourceType	Type of resource (LINK, ARTICLE, or VIDEO)

Duplicates were removed based on CourseName + TopicName + ResourceURL.

## Data Cleaning

### usm\_courses

- Compared the fields in the “courses” table, to the existing fields in “usm\_courses”.
- Created a new data frame with fields that matched that of the “courses” table, so that the data from “usm\_courses” could fit in it.
  - Fields: course\_code, course\_name, college\_id
  - The field for the primary key (course\_id) is unnecessary, since the field is auto-incrementing, and will be added automatically later via SQL indexing.

#### 1. Algorithm For Cleaning

- For each row (ie, usm course) in the CSV file, I did the following:
  - We wanted the course abbreviation concatenated with the course catalog number to be the “course code”.
  - Since every course name had the course abbreviation, and the catalog number as the first 7 characters, followed by 2 more unnecessary characters, I needed to remove those first 9 characters from every course name.
    - Before I did so, I first took the first 7 characters, and excluded the 4th character, which was white space, and saved that as the row’s course code.
  - I then took every character in the course name from the 10th character onwards, and saved that as the row’s new course name.

## 2. Saving Results

- At the end of each iteration in the for loop, I added the following values to a new row in the new courses dataframe:
  - The value of course\_code was the course code previously saved earlier in the loop.
  - The value of course\_name was the new course name previously saved earlier in the loop.
  - The value of college\_id was set to null, since the college\_id field in the schema was a nullable field.

## us\_colleges

- Mostly a repeat of the data cleaning process for “usm\_courses”, minus the different fields.
- Compared the fields in the “colleges” table, to the existing fields in “us\_colleges”.
- Created a new data frame with fields that matched that of the “courses” table, so that the data from “us\_colleges” could fit in it.
  - Fields: course\_code, course\_name, college\_id
  - The field for the primary key (college\_id) is unnecessary, since the field is auto-incrementing, and will be added automatically later via SQL indexing.

## 1. Algorithm For Cleaning

- For each row (ie, a US college) in the CSV file, I did the following:
  - Take the name of the college.
  - If the name of the college had any commas in it, replace it with an empty character.
    - This was to prevent truncation errors later on during indexing.
  - Save the new college name with no commas.

## 2. Saving Results

- At the end of each iteration in the for loop, I added the following values to a new row in the new colleges dataframe:
  - The value of college\_name was the name of the college, updated to remove any commas.

## quiz\_data

- Compared the fields in the “question” table, and the “answer” table to the existing fields in “quiz\_data”



- Created two new data frames, one for question, and one for answer, both with the same fields as their respective tables (ie,
  - Question: question\_id, quiz\_id, question\_text, question\_type, points
  - Answer: answer\_id, question\_id, answer\_text, is\_correct

## 1. Questions Data Frame

### Algorithm For Cleaning / Saving Results

- First dealt with the **questions** data frame.
- I used a for loop to loop through all of the rows of the raw quiz data, and created a new row in the questions data frame to save all of the new values for each iteration.
  - The value of question\_id is auto incrementing, so the value would just be the iteration i in the for loop.
  - quiz\_id is 0 by default, since all these questions are part of the same quiz.
  - question\_text is the same as the “question” field in “quiz\_data”.
  - The question\_type for all questions is “multiple\_choice”, since the given data only includes multiple choice questions.
  - “points” has a value of 1 by default.

## 2. Answers Data Frame

### Algorithm For Cleaning

- Next, dealt with the **answers** data frame.
  - Since the answers were multiple choice, I needed a way to separate, and clean all the individual answers in the “Answer” field of “quiz\_data”. I created a helper method to do this.
    - The entire answer field was surrounded by square brackets, so I needed to remove those.
    - Every individual answer was surrounded by an apostrophe, followed by a comma, so I would also need to remove those.
    - To do this, I created a list called “answerCharacters” that would store only the characters of the actual answer, which would then be joined into a single string.
    - To exclude square brackets, I simply made the for loop continue with the loop if it encountered any.
    - To exclude commas, it would only store the characters in the answer that were non-commas.
    - To exclude apostrophes, and to make sure every individual answer was separated, I made a variable keeping track of the apostrophe count. If the count ever reached 2, we know that we already have the full answer

stored. It would then join all the characters in “answerCharacters”, and append it to the list of answers.

- For subsequent iterations, I made sure to clear the apostrophe count, and the “answerCharacters” list.
- Finally, this method returns a list of all the different answers.
- I created a nested for loop to index all the different answers. The outer for loop is to iterate through a range equal to the number of questions, and the inner for loop is to iterate through all the answers of the question.
  - This way, I’m able to store all the answers as individual rows, and create auto-incrementing unique ids for each answer, while at the same time making sure that all of the answers’ question\_id match up with their respective question’s id.

### 3. Saving Results

- Finally, after the end of each iteration, I add a new row at the end of the new answers dataframe, and add the values for each field in the row.
  - The value of answer\_id was the current iteration of the inner for loop.
  - The value of question\_id was the current iteration of the outer for loop.
  - The value of answer\_text was the value of the current iteration of the inner for loop (ie, the current answer’s text).
  - is\_correct was set to 0 (ie, false) by default.
    - While this marks every answer as incorrect, there wasn’t really a way to automatically, and correctly mark each question as right or wrong. At least, not without using AI, which was beyond the scope of the project.
    - The reason I set everything to false instead of true, is that every multiple choice question presumably only had one correct answer, so at least this way, 3/4ths of the answers are correctly marked as false.

## course\_resources

### 1. Objectives and Justification

- This cleaner takes a raw CSV of scraped data of course-related links and turns it into:
  - A cleaned CSV
  - A MySQL Insert script to populate the Resource table
- The goal of this is to:
  - Normalize Messy URLs
    - Base64 encoded links
    - Bing redirect wrappers
  - Generate readable titles
  - Place the resource under the correct filetype

- Guarantee that every inserted row has consistent, valid data that matches the schema

## 2. Methodology & Tools

- Pandas
  - Loading CSVs, building cleaned DataFrames, and exporting a cleaned CSV
- os
- re
  - Detects base64 blobs and patterns in URLs/titles
- base64
  - Decodes base64 urls
- Urllib.parse
  - Parsing and unwrapping redirects URLs
- The python program expects an input CSV that went through the resource scraping process. The CSV should include:
  - A query or original search string used to find the resource
  - A raw title or label
  - A raw URL
- CSV is loaded with:
  - **df = pd.read\_csv(INPUT\_CSV, dtype=str, keep\_default\_na=False)**
    - Ensures all empty fields remain empty
    - All values are treated as strings

## 3. Cleaning Algorithm

- The cleaning logic has three main tasks:
  - Make URLs usable
  - Generate readable titles
  - Infer file types
- strip\_weird(s: str)
  - Replaces non-breaking spaces such as `\u00a0` and special characters like `">"` with regular spaces
  - Strips leading and trailing white space
  - Multiple spaces in a row are converted to a single space
- URL Normalization
  - Raw CSV contained URLs that were wrapped in Bing redirect links
  - Instead of http or https URLs, there were many weird base64 blobs
  - Maybe\_b64\_url(s: str)
    - If string starts with http or https, it remains the same
    - If the string does not, it searches for base64 substrings
    - Converts `"-"` to `"+"` and `"_"` to `"/"` to make substring base64 compatible
    - Looks for `"aHR0"` which is base64 for http

- unwrap\_bing(u: str)
  - Parses URL
  - If a host ends with bing.com and starts with /ck/
    - Pulls the URL from the u query
    - Unquotes it
    - Passes through maybe\_b64\_url
- normalize\_url(u: str)
  - Calls the strip\_weird function
  - If URL contains bing, calls unwrap\_bing
  - It calls maybe\_b64\_url
  - Returns a clean url
- Title Normalization
  - Many scraped titles were not understandable or did not make sense
    - Examples of this include “home” and “login” or domains like “default”
  - title\_from\_url(u: str)
    - Takes a basename or the URL path
    - Strips the extension from the URL
    - Replaces “-” and “\_” with spaces
  - clean\_title(raw\_title: str, url: str)
    - Calls strip\_weird on the title
    - Removes leading http or https
    - A set of titles that should not be used as titles are defined
      - If the cleaned title is empty, less than 3 chars, just a domain or in this set of titles it uses url\_fallback
  - url\_ballback
    - Uses a non-empty URL path if it is available and follows a <host>/<segment> template
    - Applies .title() with length limits
      - Examples include:
        - www.britannica.com / anthropology
        - mexico.internationaltrucks.com / camiones de carga
- Filetype inference
  - filetype\_from\_url(u: str)
    - Parses the URL path and extracts the file extension
    - Uses a map of known extensions
    - If the extension isn't found on the map, it defaults to “LINK”

## 4. Cleaned DataFrame

- After the data has been cleaned, the python program will build a new DataFrame
- Cdf = pd.DataFrame
  - Description is left blank for now (looking to improve upon this)
  - Title, filetype, and source are ready to load
- Drops any row where both the title and source are empty

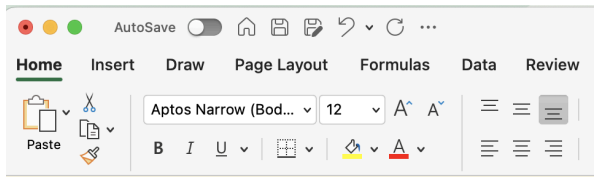
- Logs row count before and after to verify how many rows were removed
- SQL Building
  - Once the DataFrame is ready, a SQL script is generated
  - Starts with a f'SET @uploader\_id := {DEFAULT\_UPLOADER\_ID}
  - For each cleaned row it generates:
    - INSERT INTO Resource (uploader\_id, title, description, filetype, source)
    - VALUES (@uploader\_id, '<title>', '<description>', '<filetype>', '<source>');
  - Handles quotes inside strings by replacing " " with " " " " so every value has its single quotes doubled
- Not a part of the Python Program, but added after is a seeder user and uploader\_id binding
  - Since Resource.uploader\_id is a foreign key:
    - FOREIGN KEY (uploader\_id) REFERENCES Users(user\_id)
      - Ensures a system user exists and binds uploader\_id to it
    - The dedicated system user is:
      - User\_id = 1001
      - Email = resources@system.local
      - Name: System Seeder

## 5. Validation and Output

- To ensure that the cleaning and SQL generation worked:
  - Verified that the resource table in study buddy matched the columns created by the python program
  - Compared raw\_count to the number of rows in course\_resources\_cleaned.csv
  - Reviewed titles and URLs to confirm they were readable and usable
  - Ran generated insert\_resources.sql against the StudyBuddy schema and confirmed successful inserts

# Appendix

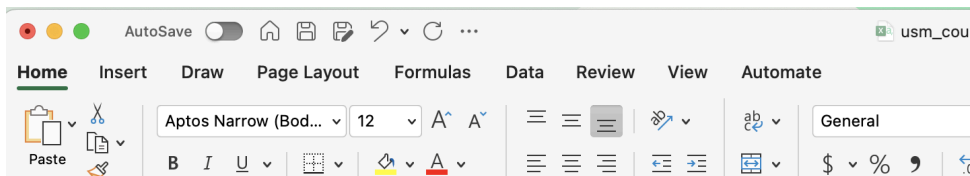
 **us\_colleges** [\[Click or Hover Link to Data\]](#)



**Possible Data Loss** Some features might be lost if you save this workbook in th

	A	B	C
1	State	College Name	
2	Connecticut	Albertus Magnus College	
3	Connecticut	Asnuntuck Community College	
4	Connecticut	Capital Community College	
5	Connecticut	Central Connecticut State University	
6	Connecticut	Charter Oak State College	
7	Connecticut	Connecticut College	
8	Connecticut	Eastern Connecticut State University	
9	Connecticut	Fairfield University	
10	Connecticut	Gateway Community College	
11	Connecticut	Goodwin University	
12	Connecticut	Hartford International University for Religion and Peace	
13	Connecticut	Holy Apostles College and Seminary	
14	Connecticut	Housatonic Community College	
15	Connecticut	Manchester Community College	
16	Connecticut	Middlesex Community College	
17	Connecticut	Mitchell College	
18	Connecticut	Naugatuck Valley Community College	
19	Connecticut	Northwestern Connecticut Community College	
20	Connecticut	Norwalk Community College	
21	Connecticut	Post University	
22	Connecticut	Quinebaug Valley Community College	
23	Connecticut	Quinnipiac University	
24	Connecticut	Sacred Heart University	

 **usm\_courses** [\[Click or Hover Link to Data\]](#)



**Possible Data Loss** Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To pres

	A	B	C	D
1	course_code	course_name	course_link	page_number
2	N/A	ANT 306-1-1Analysis of Archaeological Materials	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
3	N/A	ANT 308-1-1Environmental Archaeology	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
4	N/A	ANT 315-1-1Ethnography: Methods, Ethics, and Practi	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
5	N/A	ANT 320-1-1Anthropology & the Museum	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
6	N/A	ANT 350-1-1International Development	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
7	N/A	ANT 360-1-1Public Archaeology	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
8	N/A	ANT 380-1-1African American Historical Archaeology	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
9	N/A	ANT 395-1-1Topics in Anthropology	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
10	N/A	ANT 410-1-1Japan: Archaeology, Environmental Histo	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
11	N/A	ANT 495-1-1Topics in Anthropology	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
12	N/A	ANT 508-1-1Environmental Archaeology	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
13	N/A	ANT 515-1-1Ethnography: Methods, Ethics, and Practi	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
14	N/A	ANT 560-1-1Public Archaeology	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
15	N/A	ARA 101-1-1Beginning Arabic I	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
16	N/A	ARA 102-1-1Beginning Arabic II	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
17	N/A	ARA 150-1-1Language Table	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
18	N/A	ARA 201-1-1Intermediate Arabic I	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2
19	N/A	ARA 202-1-1Intermediate Arabic II	https://catalog.usm.maine.edu/preview_course_nopop.php?i	2

## course\_resources\_scraped [Click or Hover Link to Data]

course_resources_scraped				
CourseName	TopicName	ResourceURL	ContentSnippet	ResourceType
Data Structures	Array	<a href="https://www.geeksforgeeks.org/dsa/array-data-structure-guide/">https://www.geeksforgeeks.org/dsa/array-data-structure-guide/</a>	An array is a collection of data items stored at consecutive memory locations.	LINK
Data Structures	Linked Lists	<a href="https://www.geeksforgeeks.org/dsa/linked-list-data-structure/">https://www.geeksforgeeks.org/dsa/linked-list-data-structure/</a>	Like arrays, Linked List is a linear data structure. Unlike arrays, Linked List is a linear data structure which follows a pointer to the next node.	LINK
Data Structures	Stack	<a href="https://www.geeksforgeeks.org/dsa/stack-data-structure/">https://www.geeksforgeeks.org/dsa/stack-data-structure/</a>	Stack is a linear data structure which follows a Last In First Out (LIFO) principle.	LINK
Data Structures	Queue	<a href="https://www.geeksforgeeks.org/dsa/queue-data-structure/">https://www.geeksforgeeks.org/dsa/queue-data-structure/</a>	Like Stack, Queue is a linear structure which follows a First In First Out (FIFO) principle.	LINK
Data Structures	Binary Tree	<a href="https://www.geeksforgeeks.org/dsa/binary-tree-data-structure/">https://www.geeksforgeeks.org/dsa/binary-tree-data-structure/</a>	Unlike Arrays, Linked Lists, Stack and queues, which are linear data structures, Binary Tree is a non-linear data structure.	LINK
Data Structures	Binary Search Tree	<a href="https://www.geeksforgeeks.org/dsa/binary-search-tree-set-1-search-and-insertion/">https://www.geeksforgeeks.org/dsa/binary-search-tree-set-1-search-and-insertion/</a>	A Binary Search Tree is a Binary Tree following the Binary Search property.	LINK
Data Structures	Heap	<a href="https://www.geeksforgeeks.org/dsa/heap-data-structure/">https://www.geeksforgeeks.org/dsa/heap-data-structure/</a>	A Heap is a special Tree-based data structure in which each node has at most two children and satisfies the heap property.	LINK
Data Structures	Hash Table Data Structure	<a href="https://www.geeksforgeeks.org/dsa/hash-table-data-structure/">https://www.geeksforgeeks.org/dsa/hash-table-data-structure/</a>	Hashing is an important Data Structure which is used to store and retrieve data efficiently.	LINK
Data Structures	Matrix	<a href="https://www.geeksforgeeks.org/dsa/matrix/">https://www.geeksforgeeks.org/dsa/matrix/</a>	A matrix represents a collection of numbers arranged in rows and columns.	LINK
Data Structures	Trie	<a href="https://www.geeksforgeeks.org/dsa/trie-insert-and-search/">https://www.geeksforgeeks.org/dsa/trie-insert-and-search/</a>	Trie is an efficient information retrieval data structure.	LINK
Algorithms in Programming	What is an Algorithm   Introduction to Algorithms	<a href="https://www.geeksforgeeks.org/dsa/introduction-to-algorithms/">https://www.geeksforgeeks.org/dsa/introduction-to-algorithms/</a>	The word Algorithm means a set of finite rules or instructions to solve a problem.	LINK
Algorithms in Programming	Algorithms and Programming Concepts (PDF)	<a href="https://www.montclair.edu/computer-science-education/wp-content/uploads/2016/08/Algorithms-and-Programming-Concepts.pdf">https://www.montclair.edu/computer-science-education/wp-content/uploads/2016/08/Algorithms-and-Programming-Concepts.pdf</a>	A comprehensive, multi-page PDF document covering various algorithms and programming concepts.	PDF
Web Applications Development	Web Development Technologies	<a href="https://www.geeksforgeeks.org/web-tech/web-technology/">https://www.geeksforgeeks.org/web-tech/web-technology/</a>	Web development refers to building, creating, and maintaining websites.	LINK
Web Applications Development	Backend Development	<a href="https://www.geeksforgeeks.org/web-tech/web-technology/#backend-development">https://www.geeksforgeeks.org/web-tech/web-technology/#backend-development</a>	In this module, we will explore the technologies and concepts used in backend development.	LINK
Web Applications Development	Frontend Development	<a href="https://www.geeksforgeeks.org/web-tech/web-technology/#frontend-development">https://www.geeksforgeeks.org/web-tech/web-technology/#frontend-development</a>	In this module, we explore the core technologies and concepts used in frontend development.	LINK

## us\_majors.csv

major_name				
Accounting and Computer Science				
Accounting and Related Services				
Advanced/Graduate Dentistry and Oral Sciences				
Aerospace, Aeronautical and Astronautical Engineering				
African Languages, Literatures, and Linguistics				
Agricultural Business and Management				
Agricultural Engineering				
Agricultural Mechanization				
Agricultural Production Operations				
Agricultural Public Services				
Agricultural and Domestic Animal Services				
Agricultural and Food Products Processing				
Agriculture				
Agriculture, Agriculture Operations, and Related Sciences				
Agriculture/Veterinary Preparatory Programs				
Air Force ROTC, Air Science and Operations				
Air Transportation				
Allied Health Diagnostic, Intervention, and Treatment Professions				
Allied Health and Medical Assisting Services				
Alternative and Complementary Medical Support Services				
Alternative and Complementary Medicine and Medical Systems				
American Indian/Native American Languages, Literatures, and Linguistics				
American Sign Language				
Animal Sciences				
Anthropology				
Anthrozoology				
Apparel and Textiles				
Applied Horticulture and Horticultural Business Services				