

Name: _____ Partners: _____

Python Activity 70: Web Scraping with Python

Python is a *powerful tool* that processes & manipulates, data; accesses data; creates beautiful things such as art, solutions, and puzzles; and expands human capabilities. But it also *communicates complex computational ideas!*

Learning Objectives

Students will be able to:

Content:

- Describe what **HTML** is
- Define some common **HTML tags**
- Predict what code using the beautifulsoup4 module will do

Process:

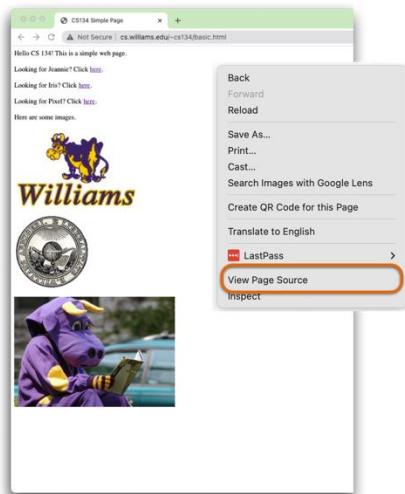
- Write some basic **HTML code**
- Write python code that scrapes web data

Prior Knowledge

- Python concepts: modules, methods

Critical Thinking Questions:

1. If we open a web page (cs.williams.edu/~cs134/basic.html) with a web browser, and then right-click the page and select "View Page Source" as below on the left, a new window opens with the text on the right (only first 11 lines shown):

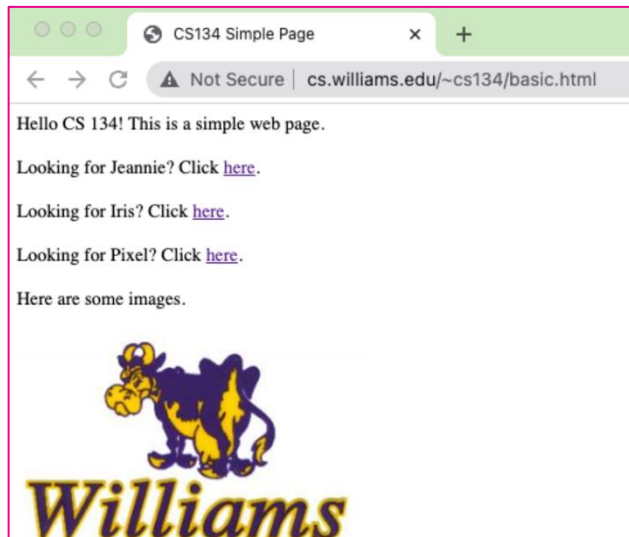


```
basic.html
<html>
<head>
  <title>CS134 Simple Page</title>
</head>
<body>
  Hello CS 134! This is a simple web page.

  <br><br>
  Looking for Jeannie? Click <a
href="http://www.cs.williams.edu/~jeannie">here</a>

  <br><br>
  Looking for Iris? Click <a
href="http://www.cs.williams.edu/~iris">here</a>.
```

- a. If we change the line "Hello CS 134! This is a simple web page." to "Hello CS 134! This is a simple web page." what do you think will happen to the webpage (zoomed-in, below)?:



- b. The text, "Hello CS 134! This is a simple web page." will turn blue. What do you think the "``" tag does?
-
- c. What do you think the "``" tag does?
-
- d. How might we modify the "``" tag to make text *green*?
-

FYI: *Hypertext Markup Language, or HTML*, is a markup language (not a programming language!) that specifies how to format text for your Internet browser. Different tags/symbols specify how the computer should display the text.

2. Match the HTML tags on the left with what you think they might do on the right:

<code><h1>Text goes here</h1></code>	Makes the text bold (or: <code></code>)
<code>Text goes here</code>	Changes the font type
<code><i>Text goes here</i></code>	Makes a hyperlink
<code>Text here</code>	Changes the text's color
<code>Text here</code>	Changes the text's size
<code>Text goes here</code>	Inserts ~2 newlines (paragraph definition)
<code>Text goes here</code>	Inserts ~1 newline (line break)
<code><p>Text goes here</p></code>	Makes a level 1 heading
<code>
</code>	Makes the text italic (or: <code></code>)

3. Well-formed HTML files also need a bit of structure. Match the structural HTML tags on the left with what you think they might do on the right:

<code><html></code>	Specifies text to appear at top of web browser
<code><head>Text & Tags go here</head></code>	HTML in here is the main part of the file
<code><title>Text goes here</title></code>	Ends the HTML file
<code><body>Text & Tags go here</body></code>	Defines what markup language is being used
<code></html></code>	HTML in here is the metadata for the file

4. We can import a new module, `requests`, and use its methods as below:

```

Interactive Python
>>> import requests
>>> r = requests.get('http://www.cs.williams.edu/~cs134/basic.html')
>>> r.text

'<html>\n <head>\n   <title>CS134 Simple Page</title>\n </head>\n\n <body>\n
Hello CS 134!  This is a simple web page.\n\n   <br><br>\n   Looking for
Jeannie?  Click <a href="http://www.cs.williams.edu/~jeannie">here</a>.\n\n
<br><br>\n   Looking for Iris?  Click <a
href="http://www.cs.williams.edu/~iris">here</a>.\n\n   <br><br>\n   Looking for
Pixel?  Click <a
href="https://www.cs.williams.edu/~iris/website/img/HAILab.jpg">here</a>.\n\n
<br><br>\n   Here are some images. \n\n   <br><br>\n   \n\n
<br><br>\n   
'

```

- a. What is `'http://www.cs.williams.edu/~cs134/basic.html'`?
Hint: How does it relate to the webpage in the previous questions?

How does the text displayed by `r.text` relate to the HTML file we saw in question 1?

- b. What do you think the `requests.get(...)` method does?

What do you think the `r.text` attribute contains?

Is what's stored in `r.text` easily readable? _____

5. We can import a new module, `beautifulsoup`, and use its methods as below:

Interactive Python (continued)

```
>>> from bs4 import BeautifulSoup
>>> soup = BeautifulSoup(r.text, 'html.parser')
>>> print(soup.prettify())

<html>
<head>
<title>
  CS134 Simple Page
</title>
</head>
<body>
Hello CS 134! This is a simple web page.
<br/>
<br/>
Looking for Jeannie? Click
<a href="http://www.cs.williams.edu/~jeannie">
  here
</a>
.
<br/>
<br/>
Looking for Iris? Click
<a href="http://www.cs.williams.edu/~iris">
```

- a. How does the text displayed by `print(soup.prettify())` relate to the `r.text` command we saw previously?

Is what's returned by `print(soup.prettify())` easily readable? _____

- b. What do you think the `BeautifulSoup` module does?
-

FYI: Oftentimes, data we obtain is not in a meaningful format. To convert raw text data into more meaningful data, we must *parse* the data. We could, of course, build our own custom string manipulation methods to implement this behavior, but there's plenty of existing Python modules that will do this for us, of which *BeautifulSoup* is one.

6. Observe these additional interactions in interactive Python that continue from our previous example:

```


Interactive Python



```
>>> soup.title
<title>CS134 Simple Page</title>
>>> soup.title.name
'title'
>>> soup.title.string
'CS134 Simple Page'
>>> soup.title.parent.name
'head'

>>> soup.img

>>> soup.a
here
>>> soup.find_all('a')
[here,
 here,
 here]
```


```

- a. What *type* of object might `soup` be?

What does the `soup.title` attribute contain?

What *type* of object might `soup.title` be?

- b. What does the `.name` attribute contain? *Hint: it does this for all HTML tags!*

What does the `.string` attribute contain? *Hint: it does this for all HTML tags!*

- c. What does the `soup.img` attribute contain?

What might `img` be an abbreviation for? _____

- d. What does the `soup.a` attribute contain?

What does the `soup.find_all('a')` method return?

7. Observe these additional interactions in interactive Python that continue from our previous example:

```
Interactive Python  
  
>>> for link in soup.find_all('a'):  
...     print(link.get("href"))  
  
http://www.cs.williams.edu/~jeannie  
http://www.cs.williams.edu/~iris  
https://www.cs.williams.edu/~iris/website/img/HAILab.jpg  
  
>>> for link in soup.find_all('a'):  
...     print(link.get_text())  
  
here  
here  
here
```

- a. What does the above code do?

8. Name at least three reasons why we might want to pull source code from web pages:

Application Questions: Use the Python Interpreter to check your work.

1. Explore the BeautifulSoup documentation: <https://beautiful-soup-4.readthedocs.io/en/latest/>
What other methods are there?

2. Write a function, `maybe_names(url, valid_words)`, that returns all invalid words found in a given webpage located at `url`. Assume that `valid_words` is a list of strings containing all valid words.