**3/9 Python Web Scraping - Grabbing a Title**

- **web scraping with Python**
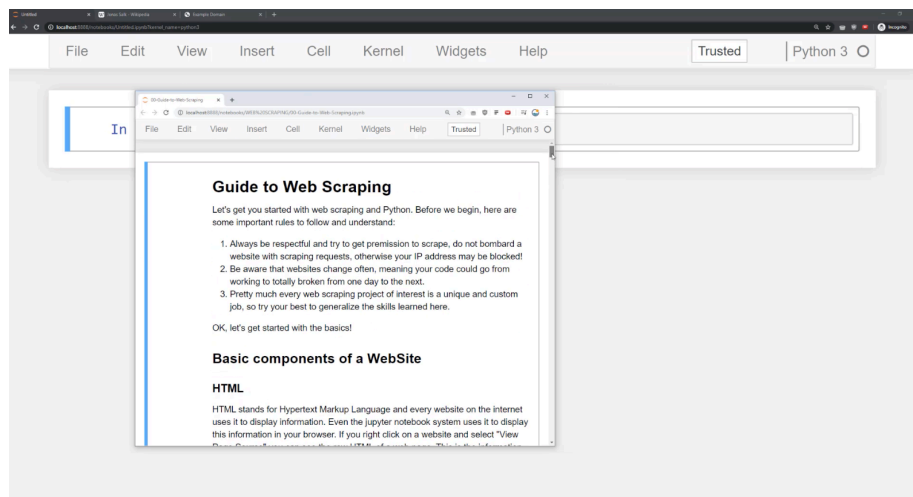  - ○ **-> we are grabbing webpage titles with Python**
    - ‣ -> in the ipynb file, it's an HTML document with web scraping
    - ‣ -> we have an example HTML webpage
    - ‣ -> he is importing the requests module
    - ‣ -> then getting a result from this
    - ‣ -> this is via the .get method
    - ‣ -> if this returns an error, then check the firewall
    - ‣ -> the request library gets a response from the URL
    - ‣ -> this has an attribute
    - ‣ -> there is information stored as a Python string
    - ‣ **-> the returns module is used to scrape information from the webpage**
    - ‣ **-> and then beautiful soup is used to format the strings**
    - ‣ -> he has then imported the bs4 module
    - ‣ -> then created a soup object, stored in the soup variable
    - ‣ -> this is passed in in xml
    - ‣ -> then calling the name of the variable returns the content which we want
    - ‣ -> we are using soup.select() to grab things from the HTML document
    - ‣ -> raw HTML elements

  - ○ **-> soup.select('')**
    - ‣ -> then passing the name of the HTML element tag into its argument
    - ‣ -> this returns a list as its default
    - ‣ -> this returns a list of all the paragraph elements on the page
    - ‣ -> we have scraped a section of a webpage, then formatted it into a more beautiful soup, and this is what we use to return a section of text from that
    - ‣ -> then we can use the .getText() method to return the string

  - ○ **-> soup.select("p"), for example**
    - ‣ -> he sets this equal to a variable

- ‣ -> he extracts the first element in the list
- ‣ **-> this is a beautiful soup object**
- ‣ -> p is the tag name in the html file
- ‣ -> by default, it returns a list of all the paragraph tags on the page
- ‣ -> for selecting html elements
- ‣ -> this also automatically returns a list
- ‣ -> we can for example add in a [0] at the end of it
- ‣ -> and then .getText()
- ‣ -> so, processing the information to extract certain terms
- ‣ -> we can also store this information in variables
- ‣ -> we also have specialised beautiful soup objects

- ○ **-> overview of steps**
  - ‣ -> he's imported requests
  - ‣ -> then requests.get
  - ‣ -> this returns a response which has a text attribute
  - ‣ -> then importing bs4
  - ‣ -> then using bs4....
  - ‣ -> we are then selecting elements off of this and passing it into a string
  - ‣ -> we need to put the right argument into .select