

From: Piruz Alemi
Subject: Web scraping & MongoDB Postings
Report on Nasa News Data Base
Date: Feb 6, 2020

Objectives met:

- Created and connected to local MongoDB databases.
- Created, read, updated, and deleted MongoDB documents using the Mongo Shell.
- Created a simple Python application that connected to and modified MongoDB databases using the PyMongo library.
- Used Beautiful Soup to scrape my own data from the web.
- Saved the results of web scraping into MongoDB.
- Became comfortable rendering templates with Flask using data retrieved from a Mongo database.
- Used Beautiful Soup to scrape data.
- Used PyMongo to save data to a Mongo database.
- Used Flask to render templates.

Helpful Links used were:

- [Mongo in 30 minutes](#)
- [Python Requests](#)
- [Webscraping with BeautifulSoup](#)
- [Python Splinter](#)

JPL Mars Space Images - Featured Image

- Visited the url for JPL Featured Space Image [here](#).
 - <https://www.jpl.nasa.gov/spaceimages/?search=&category=Mars>
- Use splinter to navigate the site and find the image url for the current Featured Mars Image and assign the url string to a variable called `featured_image_url`.
- Make sure to find the image url to the full size .jpg image.
- Make sure to save a complete url string for this image.

```
# Example:  
featured_image_url =  
'https://www.jpl.nasa.gov/spaceimages/images/largesize/PIA16225_hires.jpg'  
'
```

Completed:

- Appended the dictionary with the image url string and the hemisphere title to a list. This list contained one dictionary for each hemisphere.

```
# Example:  
hemisphere_image_urls = [  
    {"title": "Valles Marineris Hemisphere", "img_url": "..."},  
    {"title": "Cerberus Hemisphere", "img_url": "..."},  
    {"title": "Schiaparelli Hemisphere", "img_url": "..."},  
    {"title": "Syrtis Major Hemisphere", "img_url": "..."},  
,
```

Steps:

I was able to scrape the Nasa web site and post to MongoDB:

*I was able to access **children** and **siblings** of Nasa web site HTML Tags for different pages...*

Navigations of these web pages were achieved by using Splinter Library for Clicks()

Snapshots of my program follows:

```
In [43]: # go deeper into the anchor href and see what is inside?
r = requests.get("https://mars.nasa.gov/news/8558/global-storms-on-mars-launch-dust-towers-into-the-sky/", params=dict(
    query=<p>
    page=1
))
print(r.text)
```

A screenshot of a Jupyter Notebook cell showing the raw HTML content of a news article from NASA's Mars website. The code uses the requests library to get the URL and prints the resulting HTML text. The HTML includes standard headers like charset and content-type, and the main content of the news article.

I scrape various <Titles>, Divisions, <Spans>, :

```
In [47]: # Examine the results, then determine element that contains sought info
# results are returned as an iterable list
divClass = soup.findAll('div', class_="brand2")
print(divClass)
[]

In [48]: # S.P.A.N is used for in-line vs. Div.
span = soup.findAll('span')
span
```

A screenshot of a Jupyter Notebook cell showing the results of scraping spans and divs. The code finds all spans and prints them. The output shows several spans with class names like "top-nav_primary_menu_item_link_text--account", "top-nav_primary_menu_item_link_text--scores", "top-nav_primary_menu_item_link_text--news", "top-nav_primary_menu_item_link_text--video", "top-nav_primary_menu_item_link_text--standings", "top-nav_primary_menu_item_link_text--stats", and "top-nav_primary_menu_item_link_text--schedule".

Out[48]: ,

 Scores
,

 News
,

 Video
,

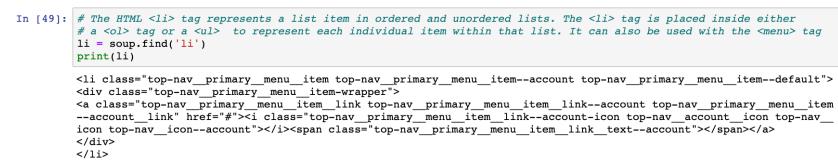
 Standings
,

 Stats
,

 Schedule
,

```
In [49]: # The HTML <li> tag represents a list item in ordered and unordered lists. The <li> tag is placed inside either
# a <ol> tag or a <ul> to represent each individual item within that list. It can also be used with the <menu> tag
li = soup.find( li )
print(li)

<li class="top-nav_primary_menu_item top-nav_primary_menu_item--account top-nav_primary_menu_item--default">
    <div class="top-nav_primary_menu_item-wrapper">
        <a class="top-nav_primary_menu_item_link top-nav_primary_menu_item_link--account top-nav_primary_menu_item--account_link" href="#">
            <i class="top-nav_primary_menu_item_link--account-icon top-nav_account_icon top-nav_top-nav_icon--account"></i>
            <span class="top-nav_primary_menu_item_link_text--account"></span>
        </a>
    </div>
</li>
```

A screenshot of a Jupyter Notebook cell showing the results of scraping list items. The code finds the first list item and prints it. The output shows an li tag with a class of "top-nav_primary_menu_item top-nav_primary_menu_item--account top-nav_primary_menu_item--default". It contains a div with a class of "top-nav_primary_menu_item-wrapper" which in turn contains an a tag with a class of "top-nav_primary_menu_item_link top-nav_primary_menu_item_link--account top-nav_primary_menu_item--account_link". This a tag has an href attribute of "#" and contains an i tag with a class of "top-nav_primary_menu_item_link--account-icon top-nav_account_icon top-nav_top-nav_icon--account". Finally, it contains a span with a class of "top-nav_primary_menu_item_link_text--account".

I was able to loop through a list of items, and post the dictionary of items:

Via Chromdrive:



On the relation between “a” and “href” :

```

In [ 1]: #### Scraped the NASA Mars News Site and collected the latest News Title and Paragraph Text.
#### Saved these parsed variables for future reference
####

In [12]: anchor_href = soup.find("a")["href"]
anchor_href

Out[12]: 'http://www.nasa.gov'

In [72]: # finding all the links
link_list = [a["href"] for a in soup.find_all('a', href=True)]
link_list

Out[72]: ['#content-wrap',
 '/',
 '/',
 'https://www.nhl.com/scores',
 '/news',
 '/news/-277384846',
 '/news/-277384234',
 '/news/-277329150',
 '/news/-278386388',
 '/news/-278387726',
 '/news/-277729393',
 '/news/-277729400',
 '/news/-277764372',
 '/news/-281011650',
 '/info/nhl/green-principles-of-principles',
 '/community/hockey-is-for-everyone',
 '/news/-277729160',
 '/news/-277549076',
 '/hockeylightsaber',
 '/info/nhl-green',
 '/info/nhl-green'
]

```

Insertions in MongoDb from a collection of Nasa News Image links and its corresponding titles follows:

_id	title	img_url
_id: ObjectId("5e3e01344ccb6408ba649dc")	Naktong Vallis	https://www.jpl.nasa.gov/spaceimages/images/wallpaper/PIA23655-640x350...
_id: ObjectId("5e3e01344ccb6408ba649dd")	Tractus Catena	https://www.jpl.nasa.gov/spaceimages/images/wallpaper/PIA23654-640x350...
_id: ObjectId("5e3e01344ccb6408ba649de")	Tarberus Fossae	https://www.jpl.nasa.gov/spaceimages/images/wallpaper/PIA23653-640x350...
_id: ObjectId("5e3e01344ccb6408ba649df")	Mangala Valles	https://www.jpl.nasa.gov/spaceimages/images/wallpaper/PIA23652-640x350...
_id: ObjectId("5e3e01344ccb6408ba649e0")	Named Valles	https://www.jpl.nasa.gov/spaceimages/images/wallpaper/PIA23651-640x350...
_id: ObjectId("5e3e01344ccb6408ba649e1")	Orcus Patera Dark Slope Streaks	https://www.jpl.nasa.gov/spaceimages/images/wallpaper/PIA23640-640x350...
_id: ObjectId("5e3e01344ccb6408ba649e2")	Barchan and Linear Dunes	https://www.jpl.nasa.gov/spaceimages/images/wallpaper/PIA23669-640x350...

But we need an array of dictionaries to insert all the data! Which we did!!
As noted in: <https://www.youtube.com/watch?v=pWbMrx5rVBE>



Next I ran the APP.py

```

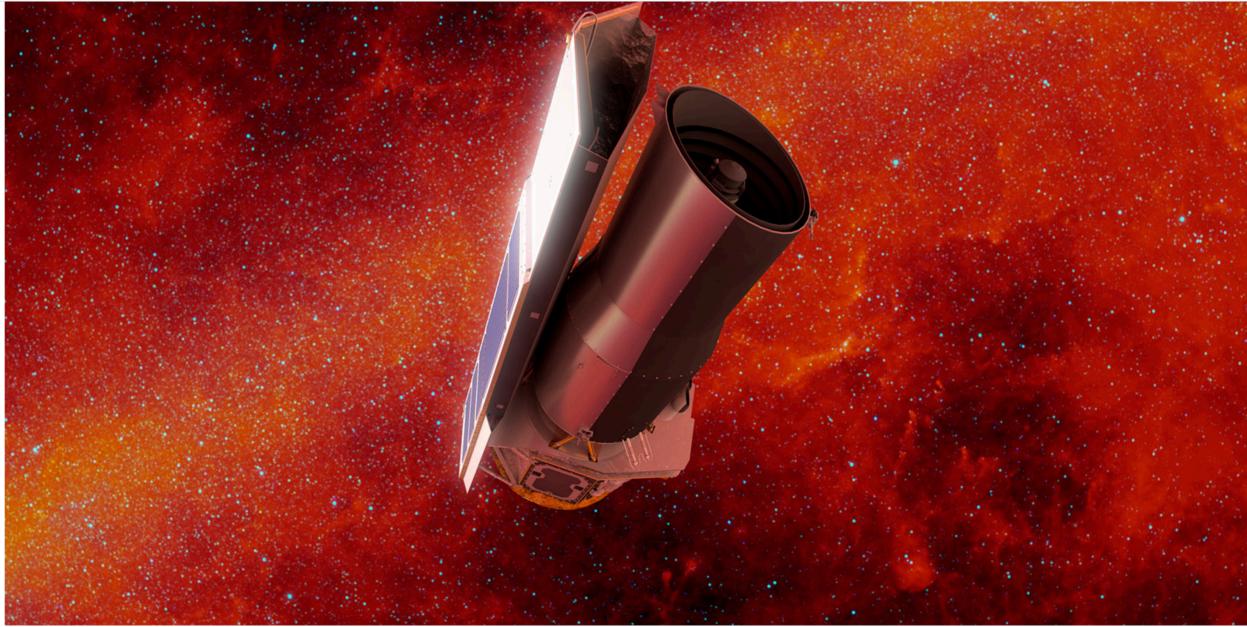
○ ○ ○ web-scraping-challenge — python • python app.py — 80x24
~/web-scraping-challenge — python • python app.py +]
app.py ~/web-scraping-challenge — python • python app.py
app2.py
app3.py
bin
export_dataframe.csv
scrapeAlemi.ipynb
scrape_costa.ipynb.py
scrape_craigslist.py
scrape_nasa.py
static
templates
~$ruz Alemi NASA Report on feb 6 2020 Images & Titles.docx
(base) Piruzs-MBP:web-scraping-challenge piruzalemi$ python app.py
* Serving Flask app "app" (lazy loading)
* Environment: production
  WARNING: This is a development server. Do not use it in a production deployment.
nt.
  Use a production WSGI server instead.
* Debug mode: on
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
* Restarting with stat
* Debugger is active!
* Debugger PIN: 153-827-043

```

Running the above APP.PY with its corresponding Template & Index-HTML under Flask and its Routes / & /Scrape
We get:



With its corresponding New Title:



The typical Mars News is:

Nasa News Title: NASA Jet Propulsion Laboratory (JPL) - Space Mission and Science News, Videos and Images

Nasa News Links: /contact_JPL.php

And a click() on Piruz Alemi Channel News, scrapes the following:

Piruz Alemi.
MARS News Channel is here for You!!! ...
Check out the latest news from Mars!
[Get MARS NEWS Data!](#)



The typical Mars News is:
Nasa News Title: NASA Jet Propulsion Laboratory (JPL) - Space Mission and Science News, Videos and Images
Nasa News Links: /contact_JPL.php