

# Learn Log - Exeter MicroApp Coding Challenge:

## Pavitra Mohandas - SVCE

### Day-1: ( Friday, 30th Oct 2020)

1. Looked into the basics of web scrapping. Read through various articles in Medium.
2. Once understanding the basics, I was able to understand the requirements better and was initially looking through how urls are extracted from a given website.
3. I came across various packages like LinkGrabber , Urllib and BeautifulSoup for performing the operation.
4. Having looked into LinkGrabber, I was able to see that along with link other attributes such as Class, id, text, seo were returned.

### Output:

```
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL

Pavitras-MacBook-Air:exeter pavitra$ python3 file.py
[ { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.kriyadocs.com/',
  'id': 'jupiter-custom',
  'seo': '',
  'text': 'Kriyadocs'},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.exeterpremedia.com/services/',
  'id': 'jupiter-custom',
  'seo': '',
  'text': ''},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.exeterpremedia.com/services/#Editorial_Services',
  'id': 'jupiter-custom',
  'seo': '#Editorial_Services',
  'text': 'EDITORIAL SERVICES'},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.exeterpremedia.com/services/#Data_Services',
  'id': 'jupiter-custom',
  'seo': '#Data_Services',
  'text': 'DATA SERVICES'},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.exeterpremedia.com/services/#Artwork',
  'id': 'jupiter-custom',
  'seo': '#Artwork',
  'text': 'ARTWORK & DESIGN'},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.exeterpremedia.com/services/#Project_Management',
  'id': 'jupiter-custom',
  'seo': '#Project_Management',
  'text': 'PROJECT MANAGEMENT'},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': '#',
  'id': 'jupiter-custom',
  'seo': '#',
  'text': '#'},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.exeterpremedia.com/blog/',
  'id': 'jupiter-custom',
  'seo': '',
  'text': 'BLOG'},
  { 'class': ['menu-item-link', 'js-smooth-scroll'],
  'href': 'https://www.exeterpremedia.com/about-us/',
  'id': 'jupiter-custom',
  'seo': ''}
```

5. Later, I looked into what internal and external links are in a website. I understood the concept and it was easy to correlate
6. Urllib's specific packages such as urlparse, urljoin for finding internal, external links.
7. I tested it with two different websites,
  1. <https://www.exeterpremedia.com/services/>
  2. <https://www.coursera.org/programs/sri-venkateswara-college-yqrql>

***Output:***

```
External link: https://www.kriyadocs.com/
Internal link: https://www.exeterpremedia.com/services/
Internal link: https://www.exeterpremedia.com/blog/
Internal link: https://www.exeterpremedia.com/about-us/
Internal link: https://www.exeterpremedia.com/events/
Internal link: https://www.exeterpremedia.com/news/
Internal link: https://www.exeterpremedia.com/contact-us/
External link: https://careers.exeterpremedia.com/jobs/Careers
Internal link: https://www.exeterpremedia.com/
Internal link: https://www.exeterpremedia.com
Internal link: https://www.exeterpremedia.com/client/
Internal link: https://www.exeterpremedia.com/exeter-privacy-policy/
Internal link: https://www.exeterpremedia.com/terms-and-conditions/
External link: tel://+91-44-23452922
External link: tel://+44-20-3287-2783
External link: tel://+1-646-736-7767
External link: mailto://sales@exeterpremedia.com
External link: https://twitter.com/Exeter_Premedia
External link: https://www.linkedin.com/company/exeter-premedia-services
External link: https://www.facebook.com/exeterpremediaservices
External link: https://www.instagram.com/exeterpremedia/
```

**Work completed Today:**

1. Getting all the links for a specific website.
2. Checking if a given URL is valid.
3. Classifying links in a given website as Internal or External.

## Day-2:(Saturday, 31st October 2020)

1. Looked on how to identify titles from a given link. Was able to get title of the given URL only.
2. Took a very long time to check how the titles can be obtained for Top Level Pages.

### *Output:*

```
Title: Exeter - Privacy Policy - Exeter
Title: Exeter - All about authors and publishers
Title: Client - Exeter
Title: Exeter - Contact Us
Title: Exeter - Events
Title: Exeter - Services
Title: Exeter - News
Title: Exeter - Blog
Title: Exeter - Happy Authors and Delighted Publishers
Title: Terms and Conditions - Exeter
```

3. Looked up online if any package was available to get homepage if a link is given as an input. Was unable to find any.
4. Derived the logic for printing the home page of a link if any other URL rather than home page is given.

### *Output:*

```
Home Page Link: https://www.exeterpremedia.com
```

5. Understood what a meta tag is. Was able to write code for getting {key, value} pairs of meta name and content.
6. The number of internal links, external links, and total was calculated at the end of the day.

### *Output:*

```
Total External links: 10
Total Internal links: 11
Total links: 21
```

7. Looked into how to obtain image urls from a given website. Discovered the package extraction. Found the image link and image size along

8. Understood the logic of website size by looking through various websites. Found the size of an entire website.

### **Word Completed Today:**

1. Page titles of all the top level pages obtained
2. Homepage given any link was derived
3. Meta Name and Meta Content values obtained
4. Total number of internal, external links obtained
5. Image URL from websites obtained
6. Website size obtained

### **Day-3:(Sunday, 1st November 2020)**

1. Looked online to check for any packages to obtain entire website data. Came to know BeautifulSoup had a function called `stripped_strings`. Was able to get data in the form of list. Converted it into string.
2. Understood what stop words are and passed the string to `stopwords` to remove stop words present in the website
3. Classified unigram and bigram from the non-stop words data.
4. Found the 20 top frequencies for unigram and bigram. Arranged them in descending order.
5. Exported values such as Home Page Link, External Link, Internal Link, Title, Total External links, Total Internal links, Website Size, Data from website, Parsed Text to text file.
6. Exported values such as Meta\_Name, Meta\_Content, Unigram, Bigram, Top 20 Unigram and Bigram, Image Url

### **Final Flow of Modules:**

User enters URL -> Validation of URL is done -> Home page of URL obtained -> All the URL's inside website is obtained -> Internal External URL's identified -> Page Titles Obtained -> Image URL obtained -> Meta Keyvalue pairs obtained -> Count of links in website -> Website size -> Data from page obtained -> Stop words removed -> With parsed text, unigram bigram obtained -> Top 20 obtained -> Data exported to CSV and Text Files

**Output:—CSV File:**

Meta Name	Meta Content	Unigram	Top 20 Frequency Unigram	Bigram	Top 20 Bigram	Image URL
viewport	width=device-width, initial-scale=1.0, minimum-scale=1.0, maximum-scale=1.0, user-scalable=0	Enter	(',', 105)	(Enter, ',')	(',', 'Wb', 15)	<a href="https://www.kootenapremedia.com/wp-content/uploads/2019/04/kootenapremedia-logo-avipg">https://www.kootenapremedia.com/wp-content/uploads/2019/04/kootenapremedia-logo-avipg</a>
format-detection	ktedge,chrome=1	-	(',', 46)	(',', 'Services')	(',', 'Our', 6)	<a href="https://www.kootenapremedia.com/wp-content/uploads/2019/04/15855784104333086-1.png">https://www.kootenapremedia.com/wp-content/uploads/2019/04/15855784104333086-1.png</a>
description	telephonesuno	Services	(Services', 15)	(Services', Kiyaboca)	(g'Editor', ',', 7)	
twittercard	Enter provides editorial services such copyediting, proofreading and indexing, XML conversion services, artwork services and project management services.	Kiyaboca	(services', 15)	(Kiyaboca', services)	(g'Editor', Services', 5)	
twitterdescription	en_US	services	(Wb', 15)	(services', EDITORIAL)	(g'Data', Services', 5)	
twittertitle	article	EDITORIAL	(',', 15)	(EDITORIAL', SERVICES)	(g'Project', Management', 5)	
twittersite	Enter - Services	SERVICES	(',', 14)	(SERVICES', DATA)	(g', ',', 5)	
twittercreator	Enter provides editorial services such copyediting, proofreading and indexing, XML conversion services, artwork services and project management services.	DATA	(',', 14)	(DATA', SERVICES)	(g'Artwork', Wb', 4)	
generator	<a href="https://www.kootenapremedia.com/services/">https://www.kootenapremedia.com/services/</a>	SERVICES	(',', 12)	(SERVICES', WORK)	(g'W', Design', 4)	
generator	Enter	ARTWORK	(Design', 12)	(ARTWORK', g')	(g', ',', 4)	
generator	summary_large_image	&	(publishers', 12)	(&', DESIGN)	(g', Author', 4)	
generator	Enter provides editorial services such copyediting, proofreading and indexing, XML conversion services, artwork services and project management services.	DESIGN	(Our', 9)	(DESIGN', PROJECT)	(g'author', publishers', 3)	
generator	Enter - Services	PROJECT	(Enter', 8)	(PROJECT', MANAGEMENT)	(g'around', time', 3)	
misapplication-TaskImage	BCenter, Premedia	MANAGEMENT	(content', 8)	(MANAGEMENT', Resource)	(g'ind', ',', 3)	
generator	BCenter, Premedia	Resources	(&', 7)	(Resource', BLOG)	(g'Wb', ',', 3)	
	Powered by LayerSlider 6.7.6 - Multi-Purpose, Responsive, Parallax, Mobile-Friendly Slider Plugin for WordPress.	BLOG	(author', 7)	(BLOG', About)	(g', 'Validating', 3)	
	WordPress 5.2	About	(booka', 7)	(About', 'a')	(g'page', 'composer', 3)	
	MasterSlider 3.2.7 - Responsive Touch Image Slider	us	(US', 6)	(us', ABOUT)	(g'Design', ',', 3)	
	User	ABOUT	(Data', 6)	(ABOUT', US)	(g'Design', services', 3)	
	May 9, 2019	US	(Artwork', 6)	(US', EVENTS)	(g'booka', 'journal', 3)	
	June 3, 2019	EVENTS		(EVENTS', NEWS)		
	Enter	NEWS		(NEWS', CONTACT)		
	Powered by WPBakery Page Builder - drag and drop page builder for WordPress.	CONTACT		(CONTACT', US)		
	Powered by Slider Revolution 5.4.8 - responsive, Mobile-Friendly Slider Plugin for WordPress with comfortable drag and drop interfaces.	US		(US', careers)		
	<a href="https://www.kootenapremedia.com/wp-content/uploads/2019/05/cropped-ktedge-1.png">https://www.kootenapremedia.com/wp-content/uploads/2019/05/cropped-ktedge-1.png</a>	careers		(careers', Kiyaboca)		
	Juniper 6.1.1	Kiyaboca		(Kiyaboca', services)		
		services		(services', EDITORIAL)		
		EDITORIAL		(EDITORIAL', SERVICES)		
		SERVICES		(SERVICES', DATA)		
		DATA		(DATA', SERVICES)		
		SERVICES		(SERVICES', WORK)		
		ARTWORK		(ARTWORK', g')		
		&		(&', DESIGN)		
		DESIGN		(DESIGN', PROJECT)		
		PROJECT		(PROJECT', MANAGEMENT)		
		MANAGEMENT		(MANAGEMENT', Resource)		
		Resources		(Resource', BLOG)		
		BLOG		(BLOG', About)		
		About		(About', 'a')		
		us		(us', ABOUT)		
		ABOUT		(ABOUT', US)		
		US		(US', EVENTS)		
		EVENTS		(EVENTS', NEWS)		
		NEWS		(NEWS', CONTACT)		
		CONTACT		(CONTACT', US)		
		US		(US', careers)		
		careers		(careers', Editor)		
		Editorial		(Editorial', Services)		
		Services		(Services', Data)		
		Data		(Data', Services)		
		Services		(Services', Artwork)		
		Artwork		(Artwork', &)		
		&		(&', Design)		
		Design		(Design', Project)		
		Project		(Project', Management)		
		Management		(Management', Editor)		
		Editorial		(Editorial', Services)		
		Services		(Services', Data)		
		Data		(Data', Services)		
		Services		(Services', Artwork)		
		Artwork		(Artwork', &)		
		&		(&', Design)		
		Design		(Design', Project)		
		Project		(Project', Management)		
		Management		(Management', Enter)		
		Enter		(Enter', ',')		
		,		(',', 'web-to-ent')		
		end-to-end		(end-to-end', Signal)		
		digital		(digital', publishing)		
		publishing		(publishing', services)		
		services		(services', designed)		
		designed		(designed', enable)		
		enable		(enable', authors)		
		authors		(authors', publishers)		
		publishers		(publishers', obtain)		

**Output — Text File:**

```

External Link: https://www.krkyaso.co.com/
Internal Link: https://www.xosterpremedia.com/servlets
Internal Link: https://www.xosterpremedia.com/blog/
Internal Link: https://www.xosterpremedia.com/about-us/
Internal Link: https://www.xosterpremedia.com/faq/
Internal Link: https://www.xosterpremedia.com/news/
Internal Link: https://www.xosterpremedia.com/contact-us/
External Link: https://careers.xosterpmedia.com/job/careers
Internal Link: https://www.xosterpmedia.com/
Internal Link: https://www.xosterpmedia.com/terms-and-conditions/
Internal Link: https://www.xosterpmedia.com/client/
Internal Link: https://www.xosterpmedia.com/xosterp-privacy-policy/
External Link: https://www.xosterpmedia.com/terms-and-conditions/
External Link: tel://491-44-2345322
External Link: tel://491-44-238-2703
External Link: tel://491-44-238-2767
External Link: mailto://@xosterpmedia.com
External Link: https://twitter.com/Xosterp_Media
External Link: https://www.linkedin.com/company/xosterp-media-services/
External Link: https://www.facebook.com/xosterpmedia/
External Link: https://www.instagram.com/xosterpmedia/

```

```
Title: Exter - Privacy Policy - Exter
Title: Exter - All about authors and publishers
Title: Client - Exter
Title: Exter - Contact Us
Title: Exter - Events
Title: Exter - Services
Title: Exter - Meet
Title: Exter - Blog
Title: Exter - Happy Authors and Delighted Publishers
Title: Terms and Conditions - Exter
```

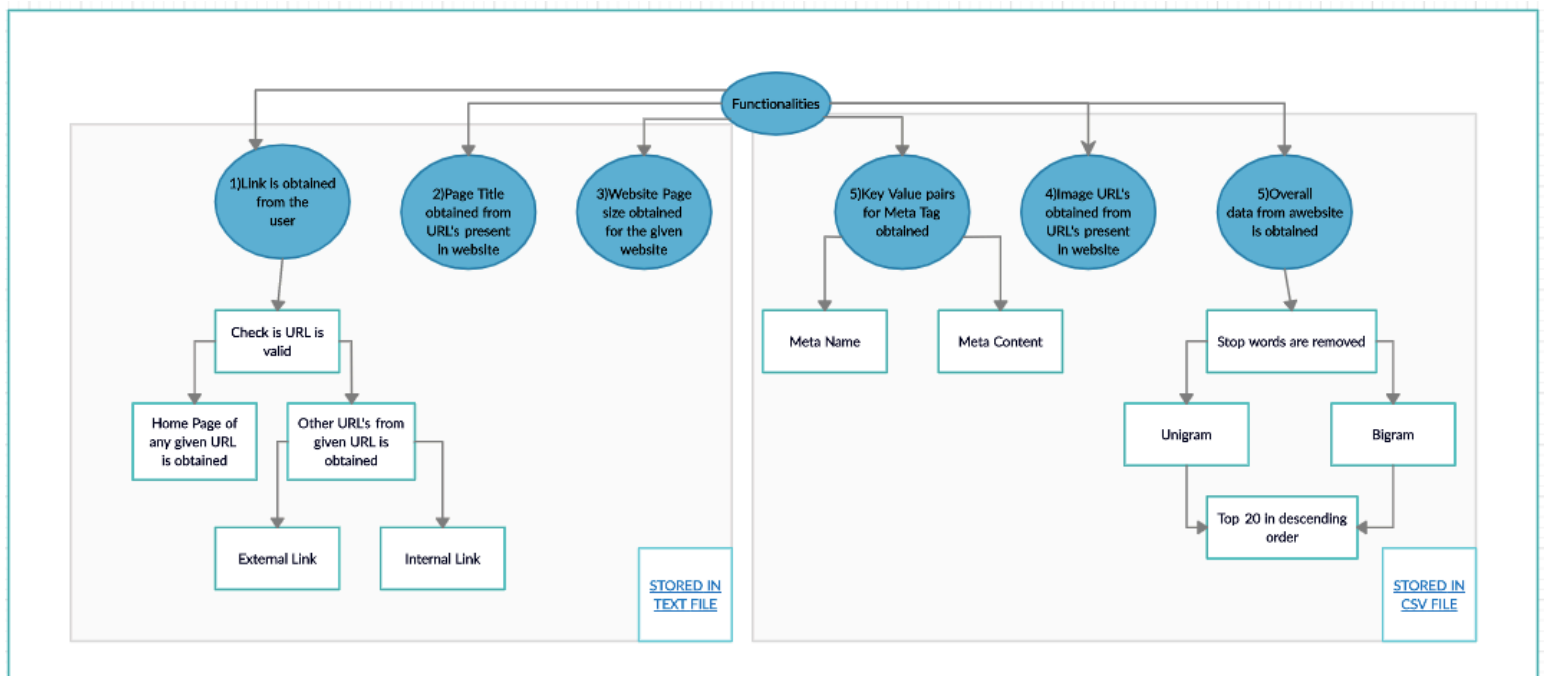
Total External links: 10  
Total Internal links: 11  
Total Links: 21

Website Size: 126.919921075 mb

[illegible][illegible]

1. Setting up virtual environment was a pretty hard task. It took sometime to setup the Virtual environment for python.
2. Every data stored into variables was of different shape. Hence I initially faced a lot of errors while writing data to the CSV file. Was able to resolve the problem when I used ignore\_text in pandas.
3. A variable called 'bi' in my code was a generator. So while passing the data to the CSV file, I was facing errors. Resolved it by converting the generator to a list.
4. Accessing title from all the url's present in a website was a very tedious process.

### Block Diagram:



In the above block diagram, I have explained all the functionalities which is available in the program that I have developed. The output file is stored in two formats based on the functionality.

### Decisions Taken:

1. Initially used linkGrabber package for extracting links. But I was unable to perform operation for finding internal and external link. So I had to scrape it and then use Urllib library for performing the task of extracting all the urls.
2. Decided to have two different files to induce more readability to the code. Establishing certain data in the form of tables made it easy to understand.
3. Made sure to have a function for checking validity of URL since there are high chances for user to enter a value which might not be a valid URL.

### **Learning:**

1. Understood the concepts related to web scrapping.
2. Learnt the packages involved in web scrapping.
3. Realized how python plays a very important role in development. It makes things really easy and the packages help save a lot of memory.
4. Various tags involved in web development.

### **Summary:**

Personally, I was very happy that I got an opportunity to be a part of the MicroApp Challenge since, I was able to effectively use 52 hours in the last 2.5 days for solving a problem. I became very curious while performing the task and got very involved into researching stuff related to Web Scrapping. I got very excited when I was getting desired outputs for the given modules. Having attended 24 hour hackathons in the past, this challenge gave me the confidence that I will be able to attend 60 hour hackathons in future. After a very long time, my boredom came to an end and I was very delighted that I was able to finish the task which was given to me.