

This is CSCCI P-14300.

# Web Scraping

what?

useful?

procedure

l: get raw html from url

2: process html

l: get raw html from url



requests

```
$ pip install requests
$ python
>>> import requests
>>> url = "https://cs50.github.io/summer/syllabus"
>>> print requests.get(url)
<Response [200]>
>>> print requests.get(url).text
<!DOCTYPE html>
...
```

2: process html

beautiful soup

```
$ pip install bs4
$ python
>>> from bs4 import BeautifulSoup
>>> url = "https://cs50.github.io/summer/syllabus"
>>> raw_html = requests.get(url).text
>>> soup = BeautifulSoup(raw_html)
```

tasks

find all links

```
>>> for link in soup.find_all('a', href=True):  
>>>     print(link.get('href'))  
/summer  
http://www.summer.harvard.edu/  
...
```



get text from paragraphs

```
>>> for para in soup.find_all('p'):
```

```
>>>     print(para.get_text())
```

```
Introduction to Web Programming
```

```
Harvard Summer School
```

```
...
```

let's get started!

```
$ git clone https://github.com/daveyproctor/web-scraping-seminar.git
```

This is CSCCI P-14300.