

Web-Scraping:
How Charlotte Is A Life Saver Again

Fayang Pan

December 3, 2013

Contents

1	An Evaluation of Web-Scraping	4
1.1	An Overview of Web-Scraping	5
1.2	Usage of Web-Scraping	5
1.3	Advantages of Web-Scraping	6
1.3.1	Speed	6
1.3.2	Convenience for Data Analysis	6
1.3.3	Availability	7
1.4	Challenges of Web-Scraping	7
1.4.1	Legality	7
1.4.2	Intricate and Dynamic Web Structures	8
1.4.3	Instability	9
2	My Internship	11
2.1	Overview of My Internship at Pinnacle Solutions, Inc.	11
2.2	How Web-Scraping Became the Best Option	12
2.3	Scrapy, a Python Framework	13
3	Scrapy Explained in Greater Detail	15
3.1	Test How Much You Can See	15
3.2	Test If the Data is Scrapable	17

3.3	Iterators	17
3.4	Specification	19
3.5	Produce a sample output	20
4	Future Development of the Project	23
4.1	Natural Language Toolkit	23
4.2	Data Mining	24
4.3	Graphic User Interface	25
	Appendices	26

List of Figures

1.1	Comparison of same url in different browsers	10
3.1	Comparison of views	16
3.2	How to extract XPath	18
3.3	Scrapy screenshot in Mac Terminal	21
3.4	Sample csv output	22

Chapter 1

An Evaluation of Web-Scraping

Introduction

Web-scraping is a way to collect existing information from the Internet. It is a nascent, powerful, but disputable subject in the arena of modern technology. Using tools to simulate viewing and downloading of a webpage through a browser, the ‘spider’ (a metaphoric name for the part that goes into each web structure) goes into the webpage and ‘crawls’ (a metaphor for copy-and-paste) desirable contents into a formatted data structure. Since this technology only requires access to the webpages, theoretically, anything that can be downloaded can be scraped off the Internet.

In the past summer, I worked with Pinnacle Solutions, Inc.¹, a company in Indianapolis, Indiana. My job was solely to develop web-scrappers. In this paper, I will provide an overview and discussion of web-scraping, followed by a description of my internship experience.

¹For information, visit <http://psiconsultants.com/>

1.1 An Overview of Web-Scraping

In general, there are three basic ways to collect data: from primary sources, e.g., conducting surveys and studies; from secondary sources, e.g., databases like the U.S. Census; and from existing but uncompiled information, e.g., product reviews from *Amazon.com*. Web-scraping focuses on the last way: to harvest data from a structured webpage, e.g., an HTML table, to a structured format, e.g., json and csv. It provides an amalgamation of the deluge of data spread over websites.

1.2 Usage of Web-Scraping

As web-scraping does not generate new data or results, it seems useless. However, people from various backgrounds may find web-scraping very useful. Here are few examples:

If a student were to conduct a study on what lessons people learn from documentaries such as *Food, Inc.*, he/she could go to <http://imdb.com>, <http://amazon.com> and <http://rottentomatoes.com> to scrape all the viewers' reviews. As there are thousands of reviews for the movie, web-scraping will help the student save much time trying to copy and paste every review. From the integrated results, the students might do a word frequency count to guess what kind of impression the majority have.

If a company were to have launched a product, they could use web-scraping to monitor their product, their competitors' products, and their partners products. If people from Amazon.com want to know how successful is the new Kindle Fire HDX, they can web-scrape all the reviews about Kindle Fire HDX, about Nexus 7, and about Kindle Fire HD. Through analysis, Amazon.com is able to collect customers' feedback promptly.

If a statistician were to monitor the price change of thousands of commodities, and if all these commodities are sold online, the statistician can use web-scraping to collect the prices. The scraper will just copy the prices from the product page of an online retailer website,

and paste all of them in a spreadsheet. Shifts in prices may reflect inflation/deflation, and the data can be used in economic research.

If a data analyst from Apple were to study how the public responds whenever a new product is released, he/she could collect the review dates from review websites and plot a frequency graph. The number of reviews generated per day may reflect how receptive the customers are to a new product.

1.3 Advantages of Web-Scraping

1.3.1 Speed

Web-scraping allows people to copy and paste from webpages faster. For instance, one wants to harvest 10,000 pieces of Amazon.com product reviews. Doing the job manually by copying and pasting every little piece of information into an excel spreadsheet would take a very long time and be prone to mistakes. Web-scraping technology, on the other hand, makes the process much faster. By specifying certain paths in the webpage structure, the spiders go into specific attributes and fields to crawl information. With the help of web-scraping, 10,000 reviews can be copied and pasted into a csv file within 15 minutes.²

1.3.2 Convenience for Data Analysis

Web-scraping improves the efficiency of data analysis. More often than not, tables of data are more ready to be processed than plain text. As web-scraping grabs from a structured format, it is able to copy and paste all the information needed into another structured format. For instance, after a csv file is generated by scraping data from Amazon product reviews, all rows in the *ratings* column have a range of 1 to 5 for number of stars, all rows in the *review*

²The result is based on running of my own scraper.

body column will contain texts of the review bodies. A well-structured database will then provide foundation for further data Data Mining and other modeling techniques.

1.3.3 Availability

Web-scraping are being more widely available today. Popular tools for web-scraping include browser extensions such as *firebug* and *iMacros*, programs such as *Wget*, programming languages such as *Perl*, *Hadoop* and *Java* which have libraries to support sending HTTP requests, and programming language extensions such as *Scrapy* and *Beautiful Soup* for *Python*. It is becoming more feasible for anyone to learn web-scraping, and to harvest data beyond what traditional ways such as web API can offer.

1.4 Challenges of Web-Scraping

1.4.1 Legality

Web-scraping raises legal and ethical issues. While it is free to download many webpages, usage of the data within might not be in compliance with the terms and conditions of the source.

There exists a fine line between using and stealing data. Currently, there is no law forbidding web-scraping all together, nor is there any way to ban downloading of webpages. However, as web-scraping becomes more versatile and powerful, companies are becoming more vigilant on data protection.

In 2000, Bidder's Edge, a website that collects auction information from various websites, collected and displayed auction information from eBay. eBay sued Bidder's Edge for "trespass to chattels"³, and Bidder's Edge soon went out of business.⁴

³For more information, visit <http://definitions.uslegal.com/t/trespass-to-chattels/>

⁴For information on the ruling, visit <http://www.tomwbell.com/NetLaw/Ch06/eBay.html>

In March 2013, Associated Press (AP) won a lawsuit against Meltwater, a Norwegian group, over copyright infringement. Before then, if the customer would like to know how frequently a certain word appears in the news, Meltwater could do that through scraping from various news websites. Meltwater used to be able to collect the news, put them into databases, and render them to customers on demand, but not any more. The court ruled that the usage of those news were beyond the scope of “fair usage”,⁵ forbidding Meltwater to scrape from news agencies since then.⁶

In general, there are few things any scraper should watch out for. First, terms of use of the websites. Even though the terms are usually long and tedious, “I did not see those” cannot be the excuse of violating the terms. Second, the purpose of the scraped information. If the scraper were to use the scraped information for malicious purposes such as spamming and libel, it is likely to be illegal. Third, possible influence to others. If the number of requests for a certain scraping activity is too large for the website server to handle, other visitors to this website may experience malfunction of the website. It is thus always wise to consult lawyers before making any scraping activity commercial.

1.4.2 Intricate and Dynamic Web Structures

As more advanced web structures emerge, scraping becomes more difficult. AJAX, for instance, sends out XMLHttpRequest only when the user performs certain maneuvers in the browser. Also, embedded JavaScript contents may not be downloaded locally. The advent of these dynamic ways of loading web contents brings new challenges to web-scraping.

⁵For more information about fair use, visit <http://fairuse.stanford.edu/overview/fair-use/four-factors/>

⁶For more information on the ruling, visit <http://www.scribd.com/doc/131847330/Meltwater-AP-Ruling>

1.4.3 Instability

Another underlying problem comes from the fact that web-scraping is based on webpage structure. If the website adds a few toolbars, the locations of HTML attributes will change, resulting in crash of the scraper. There also exists the problem of encoding and decoding. If the website has an unusual encoding for text, the scraper may crash, too. If the structure is dependent on web browser, they content might be different as well. These potential problems manifest web-scraping's intrinsic dependence on web interface.

In the following two pictures, the same webpage is opened using two different browsers, Google Chrome and Mozilla Firefox, under the same environment at the same time. They are different in locations of choice of flavor, availability of customer image, font sizes and layouts, sites available for sharing, etc.

amazon

[Try Prime](#)
[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

Shop by Department

Search

Grocery & Gourmet Food

Go

[Hello, Sign in Your Account](#)
[Try Prime](#)
[Cart](#)
[Wish List](#)

[Grocery & Gourmet Food](#)
[Best Sellers](#)
[Breakfast Foods](#)
[Beverages & Coffee](#)
[Snack Food](#)
[Baking](#)
[Edible Gifts](#)
[Fresh Flowers](#)
[Candy](#)
[Special Offers](#)
[Subscribe & Save](#)

Lindt Lindor Truffles, 60 Count

by Lindt

★★★★★
(376 customer reviews)

Flavor: Milk

Product Packaging: Standard Pack

Price: **\$17.99** (\$0.30 / count) & **FREE Shipping** on orders over \$35. [Details](#)

In Stock.
 Sold by [Special Supply](#) and [Fulfilled by Amazon](#).

Want it Wednesday, Nov. 20? Order within **16 hrs 21 mins** and choose **One-Day Shipping** at checkout. [Details](#)

- Milk chocolate with a smooth filling
- Great gift for milk chocolate lovers for any occasion
- 60 count

6 new from \$16.99

Qty: 1

☐ Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

Add to Cart

or

[Sign in](#) to turn on 1-Click ordering.

Add to Wish List

More buying choices

General Store Depot	Add to Cart
\$16.99 + \$8.49 shipping (\$0.28 / count)	
SweetDeals4U LLC	Add to Cart
\$20.99 + \$7.56 shipping (\$0.35 / count)	
Omega Galleries	Add to Cart
\$24.80 + \$3.99 shipping (\$0.41 / count)	

6 new from \$16.99

[Share](#)
[Email](#)
[Facebook](#)
[Twitter](#)
[Pinterest](#)

Frequently Bought Together

amazon

[Try Prime](#)
[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

Shop by Department

Search

Grocery & Gourmet Food

Go

[Hello, Sign in Your Account](#)
[Try Prime](#)
[Cart](#)
[Wish List](#)

[Grocery & Gourmet Food](#)
[Best Sellers](#)
[Breakfast Foods](#)
[Beverages & Coffee](#)
[Snack Food](#)
[Baking](#)
[Edible Gifts](#)
[Fresh Flowers](#)
[Candy](#)
[Special Offers](#)
[Subscribe & Save](#)

Lindt Lindor Truffles, 60 Count

by Lindt

★★★★★
(376 customer reviews)

Price: **\$17.99** (\$0.30 / count) & **FREE Shipping** on orders over \$35. [Details](#)

In Stock.
 Sold by [Special Supply](#) and [Fulfilled by Amazon](#). Gift-wrap available.

Want it Wednesday, Nov. 20? Order within **16 hrs 20 mins** and choose **One-Day Shipping** at checkout. [Details](#)

Flavor: Milk

Product Packaging: Standard Pack

- Milk chocolate with a smooth filling
- Great gift for milk chocolate lovers for any occasion
- 60 count

6 new from \$16.99

[Share](#)
[Email](#)
[Facebook](#)
[Twitter](#)

Quantity: 1

☐ Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

Add to Cart

or

[Sign in](#) to turn on 1-Click ordering.

Add to Wish List

More Buying Choices

General Store Depot	Add to Cart
\$16.99 + \$8.49 shipping (\$0.28 / count)	
SweetDeals4U LLC	Add to Cart
\$20.99 + \$7.56 shipping (\$0.35 / count)	
Omega Galleries	Add to Cart
\$24.80 + \$3.99 shipping (\$0.41 / count)	

6 new from \$16.99

Frequently Bought Together

amazon

[Try Prime](#)
[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

Shop by Department

Search

Grocery & Gourmet Food

Go

[Hello, Sign in Your Account](#)
[Try Prime](#)
[Cart](#)
[Wish List](#)

[Grocery & Gourmet Food](#)
[Best Sellers](#)
[Breakfast Foods](#)
[Beverages & Coffee](#)
[Snack Food](#)
[Baking](#)
[Edible Gifts](#)
[Fresh Flowers](#)
[Candy](#)
[Special Offers](#)
[Subscribe & Save](#)

Lindt Lindor Truffles, 60 Count

by Lindt

★★★★★
(376 customer reviews)

Price: **\$17.99** (\$0.30 / count) & **FREE Shipping** on orders over \$35. [Details](#)

In Stock.
 Sold by [Special Supply](#) and [Fulfilled by Amazon](#). Gift-wrap available.

Want it Wednesday, Nov. 20? Order within **16 hrs 20 mins** and choose **One-Day Shipping** at checkout. [Details](#)

Flavor: Milk

Product Packaging: Standard Pack

- Milk chocolate with a smooth filling
- Great gift for milk chocolate lovers for any occasion
- 60 count

6 new from \$16.99

[Share](#)
[Email](#)
[Facebook](#)
[Twitter](#)

Quantity: 1

☐ Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

Add to Cart

or

[Sign in](#) to turn on 1-Click ordering.

Add to Wish List

More Buying Choices

General Store Depot	Add to Cart
\$16.99 + \$8.49 shipping (\$0.28 / count)	
SweetDeals4U LLC	Add to Cart
\$20.99 + \$7.56 shipping (\$0.35 / count)	
Omega Galleries	Add to Cart
\$24.80 + \$3.99 shipping (\$0.41 / count)	

6 new from \$16.99

Frequently Bought Together

amazon

[Try Prime](#)
[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

Shop by Department

Search

Grocery & Gourmet Food

Go

[Hello, Sign in Your Account](#)
[Try Prime](#)
[Cart](#)
[Wish List](#)

[Grocery & Gourmet Food](#)
[Best Sellers](#)
[Breakfast Foods](#)
[Beverages & Coffee](#)
[Snack Food](#)
[Baking](#)
[Edible Gifts](#)
[Fresh Flowers](#)
[Candy](#)
[Special Offers](#)
[Subscribe & Save](#)

Lindt Lindor Truffles, 60 Count

by Lindt

★★★★★
(376 customer reviews)

Price: **\$17.99** (\$0.30 / count) & **FREE Shipping** on orders over \$35. [Details](#)

In Stock.
 Sold by [Special Supply](#) and [Fulfilled by Amazon](#). Gift-wrap available.

Want it Wednesday, Nov. 20? Order within **16 hrs 20 mins** and choose **One-Day Shipping** at checkout. [Details](#)

Flavor: Milk

Product Packaging: Standard Pack

- Milk chocolate with a smooth filling
- Great gift for milk chocolate lovers for any occasion
- 60 count

6 new from \$16.99

[Share](#)
[Email](#)
[Facebook](#)
[Twitter](#)

Quantity: 1

☐ Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

Add to Cart

or

[Sign in](#) to turn on 1-Click ordering.

Add to Wish List

More Buying Choices

General Store Depot	Add to Cart
\$16.99 + \$8.49 shipping (\$0.28 / count)	
SweetDeals4U LLC	Add to Cart
\$20.99 + \$7.56 shipping (\$0.35 / count)	
Omega Galleries	Add to Cart
\$24.80 + \$3.99 shipping (\$0.41 / count)	

6 new from \$16.99

Frequently Bought Together

amazon

[Try Prime](#)
[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

Shop by Department

Search

Grocery & Gourmet Food

Go

[Hello, Sign in Your Account](#)
[Try Prime](#)
[Cart](#)
[Wish List](#)

[Grocery & Gourmet Food](#)
[Best Sellers](#)
[Breakfast Foods](#)
[Beverages & Coffee](#)
[Snack Food](#)
[Baking](#)
[Edible Gifts](#)
[Fresh Flowers](#)
[Candy](#)
[Special Offers](#)
[Subscribe & Save](#)

Lindt Lindor Truffles, 60 Count

by Lindt

★★★★★
(376 customer reviews)

Price: **\$17.99** (\$0.30 / count) & **FREE Shipping** on orders over \$35. [Details](#)

In Stock.
 Sold by [Special Supply](#) and [Fulfilled by Amazon](#). Gift-wrap available.

Want it Wednesday, Nov. 20? Order within **16 hrs 20 mins** and choose **One-Day Shipping** at checkout. [Details](#)

Flavor: Milk

Product Packaging: Standard Pack

- Milk chocolate with a smooth filling
- Great gift for milk chocolate lovers for any occasion
- 60 count

6 new from \$16.99

[Share](#)
[Email](#)
[Facebook](#)
[Twitter](#)

Quantity: 1

☐ Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

Add to Cart

or

[Sign in](#) to turn on 1-Click ordering.

Add to Wish List

More Buying Choices

General Store Depot	Add to Cart
\$16.99 + \$8.49 shipping (\$0.28 / count)	
SweetDeals4U LLC	Add to Cart
\$20.99 + \$7.56 shipping (\$0.35 / count)	
Omega Galleries	Add to Cart
\$24.80 + \$3.99 shipping (\$0.41 / count)	

6 new from \$16.99

Frequently Bought Together

amazon

[Try Prime](#)
[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

Shop by Department

Search

Grocery & Gourmet Food

Go

[Hello, Sign in Your Account](#)
[Try Prime](#)
[Cart](#)
[Wish List](#)

[Grocery & Gourmet Food](#)
[Best Sellers](#)
[Breakfast Foods](#)
[Beverages & Coffee](#)
[Snack Food](#)
[Baking](#)
[Edible Gifts](#)
[Fresh Flowers](#)
[Candy](#)
[Special Offers](#)
[Subscribe & Save](#)

Lindt Lindor Truffles, 60 Count

by Lindt

★★★★★
(376 customer reviews)

Price: **\$17.99** (\$0.30 / count) & **FREE Shipping** on orders over \$35. [Details](#)

In Stock.
 Sold by [Special Supply](#) and [Fulfilled by Amazon](#). Gift-wrap available.

Want it Wednesday, Nov. 20? Order within **16 hrs 20 mins** and choose **One-Day Shipping** at checkout. [Details](#)

Flavor: Milk

Product Packaging: Standard Pack

- Milk chocolate with a smooth filling
- Great gift for milk chocolate lovers for any occasion
- 60 count

6 new from \$16.99

[Share](#)
[Email](#)
[Facebook](#)
[Twitter](#)

Quantity: 1

☐ Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

Add to Cart

or

[Sign in](#) to turn on 1-Click ordering.

Add to Wish List

More Buying Choices

General Store Depot	Add to Cart
\$16.99 + \$8.49 shipping (\$0.28 / count)	
SweetDeals4U LLC	Add to Cart
\$20.99 + \$7.56 shipping (\$0.35 / count)	
Omega Galleries	Add to Cart
\$24.80 + \$3.99 shipping (\$0.41 / count)	

6 new from \$16.99

Figure 1.1: What Google Chrome displays versus what firefox displays with identical urls.

10

Chapter 2

My Internship

2.1 Overview of My Internship at Pinnacle Solutions, Inc.

Pinnacle Solutions, Inc.(PSI) is a data company which helps clients to analyze and interpret their databases. Founded in 1996 by Kalamazoo College alumnus Donald Penix, Jr., the company, located in downtown Indianapolis, Indiana, currently has over 20 employees. PSI primarily uses *SAS* for processing data, and *Amazon Web Services* for server management.

I was a summer intern at PSI in 2011, and my job title was Business Intelligence Content Developer. At that time, Mr. Penix was very interested in sentiment analysis, a technology aims to discover and analyze sentiment reflected through comments, reviews, and blog posts alike. For instance, a comment such as “I love this iPad!” indicates a positive sentiment, and one such as “I am sick of the iPad!” suggests a negative sentiment. Sentiment analysis will consolidate related comments and analyze them using Data Mining and NLP. After research, we found sentiment analysis a tough job for PSI, as none of the employees had the needed expertise. So we decided to start from scratch, and data collection is the first step towards sentiment analysis.

2.2 How Web-Scraping Became the Best Option

Analysis of data starts from collection of them. As previously discussed, there are three ways to collect data. PSI does not have the resources to conduct primary research, nor does the company has databases of information for reference, so it becomes evident that PSI has to collect data from existing but uncompiled data from the Internet. In order to harvest usable data from the Internet, there are many ways other than web-scraping.

First, one can extract necessary information through web API. Many websites offer APIs for developers to use to build applications. Twitter, for instance, published a new version of API months ago. Through an API, users can request information directly from the website's databases, and the data will be clean and formatted. (Twitter, for instance, offers exportable data in json format.) However, to obtain more data, users may need to pay for them. Moreover, many websites may have limited information provided in their web APIs, or may not provide APIs at all. In that case, even though users can view the data from the browser, they cannot download it through APIs. Therefore, for a small company, it might be expensive and insufficient to rely on APIs for data harvesting.

Second, there are available third-party scraping/data harvesting services. There are applications like *Mozenda*¹, and highly customizable services such as *scrapinghub*², *GNIP*³, and *Open Amplify*⁴. However, these services are not cheap to obtain. Moreover, in the long run, as PSI needs more information from the Internet, the dependency on these services will be ever increasing. Therefore, as a data company, it would be a wise choice if PSI could come up with a low cost, independent system to fetch information that is so close yet so far from them.

Third, as they are a certified *SAS* reseller, they could use the software *SAS Sentiment*

¹<http://mozenda.com/>

²<http://scrapinghub.com/>

³<http://gnip.com/>

⁴<http://www.openamplify.com/>

*Analysis Studio*⁵. However, this application is not only expensive for clients to buy, but also dependent on the web API.

After many thorough discussion sessions, a free and open software capable of fetching data from the Internet was determined to be desirable, and web-scraping seemed most promising. Many web-scraping resources are free, so the cost of doing web-scraping is low. As many programming languages and their extensions are open sources, many people contribute to the library and trouble-shooting, so the technical support is very active.

After some research, we chose Scrapy, a Python web-scraping framework, for the job.

2.3 Scrapy, a Python Framework

To quote from its website, “Scrapy is a fast high-level screen scraping and web crawling framework, used to crawl websites and extract structured data from their pages. It can be used for a wide range of purposes, from data mining to monitoring and automated testing.”⁶

Scrapy depends not only on Python, but also few other libraries. Scrapy uses Twisted⁷, an event-driven networking engine, and Zope⁸, a web application server framework. The installation guide can be found on Scrapy’s website⁹.

There are few basic ideas or classes in the Scrapy framework.

First, spider. *Spider.py* is responsible for downloading webpages, going into them and crawling data from them. At the very least, the user needs to specify the domain, the starting url, and the location in the webpage from which to crawl. If the user wants to scrape from multiple tables or multiple pages, or even conditionally scrape from certain pages, he/she can define those in *Spider.py*.

⁵For more information, visit <http://www.sas.com/text-analytics/sentiment-analysis/>

⁶<http://scrapy.org/>

⁷<http://twistedmatrix.com/trac/>

⁸<http://zope2.zope.org/about-zope-2/what-is-zope-2>

⁹<http://doc.scrapy.org/en/latest/intro/install.html>

Second, items. *Items.py* is a class specifying the column names in the result. In the end, the result would be a table consisting of all the scraped data, and items defines what are they. Items in the same class will be in the same table.

Third, settings. *Settings.py* includes information such as what information will be written into the log file, limit on frequency of requests sent, and other customizable features. In the event of a user login is needed, the user can pre-fill that in the login to avoid access denial.

Fourth, pipeline. *Pipelines.py* decides how the scraped data are processed. By default, all data will be dumped into a single file. However, if the user wants to add a filter to get rid of some data, pipeline handles that. If the user wants to split the results from one spider into two separate files, he/she can do that through coding in *pipelines.py*.

Scrapy has more advanced features, which would fall outside of the scope of this paper. In the next chapter, I would like to offer a high level approach of web-scraping using Scrapy.

Chapter 3

Scrapy Explained in Greater Detail

3.1 Test How Much You Can See

Web-scraping depends on whatever the browser can "see", and there are times when the browser sees less than a person does. To check, one can use the following line of code in the command-line (cmd.exe in Windows, Terminal in Mac OS/Linux):

```
scrapy shell http://amazon.com
```

After Scrapy finish analyzing the web address, it goes into its prompt and asks for further instructions. Here, the user can key in

```
view(response)
```

A browser window opens and shows you what the website is like in Scrapy's perspective. The following two screenshots demonstrate the differences in downloading the page through Scrapy and viewing the page through the browser.

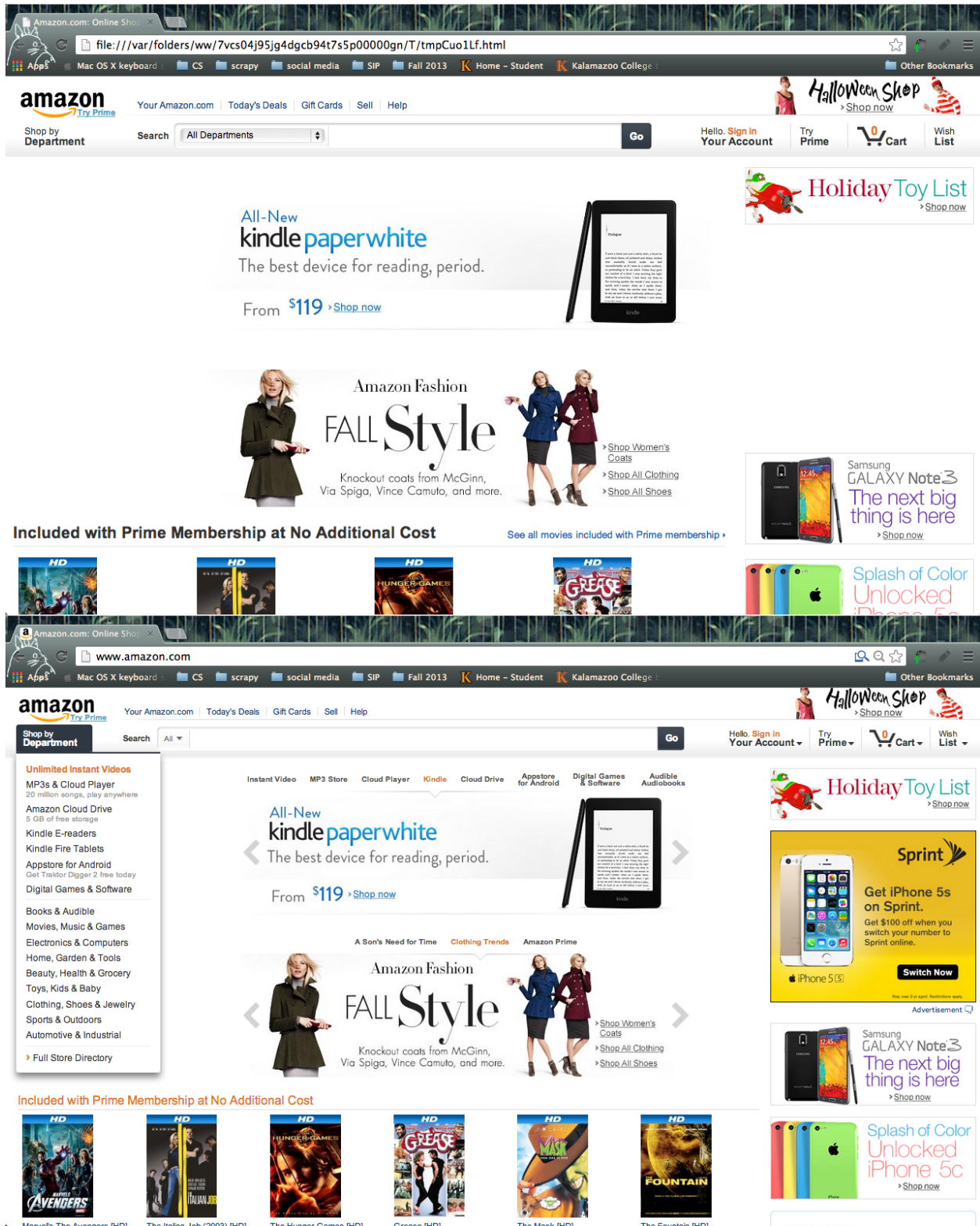


Figure 3.1: What Scrapy is able to download locally versus what the user can see.

3.2 Test If the Data is Scrapable

While in theory, any data which can be seen can be downloaded, and any downloaded data can be scraped, in practice, it is not that ideal. Even if Scrapy can see the data the user needs, the difficulty of scraping them varies. For instance, in the event of embedded JavaScript code, getting the result of a function call will be more difficult.

After setting up a project with few commands¹, it is wise to test if certain data are scrapable. Before knowing that, one needs to specify where exactly the data are located. To establish a standard way of referring to a certain field in the webpage, XPath² is used.

To retrieve a certain XPath, one can use an built-in “*Inspect Element*” tool, or an extension such as *firebug*³, available for *Firefox* and *Google Chrome*. Here is an example using “*Inspect Element*” in *Google Chrome* to extract an XPath from `http://www.nbcnews.com`:

After knowing a specific XPath, one can write related code to test if Scrapy can recognize the XPath and scrape the data. If one piece of data is scrapable, it is likely that the entire table which contains that piece of data can be scraped. If the table is scrapable, it is also likely that similar tables can be scraped, too.

3.3 Iterators

Web-scraping relies heavily on the structure of the webpage. If the webpage is very organized, it is possible for the user to use XPath to find all entries in a table, and send out a request to scrape the next page. In Scrapy, the request is in the form of a callback function. The spider will scrape the data yield the results from the page, and then execute the callback function. Here is a snippet of that concept:

¹For more information, visit <http://doc.scrapy.org/en/latest/intro/tutorial.html>. The website has a simple example on how to start a project

²Abbreviation of “XML Path Language”. For more information, see <http://www.w3.org/TR/xpath/>

³Visit <http://getfirebug.com/formoreinformation>

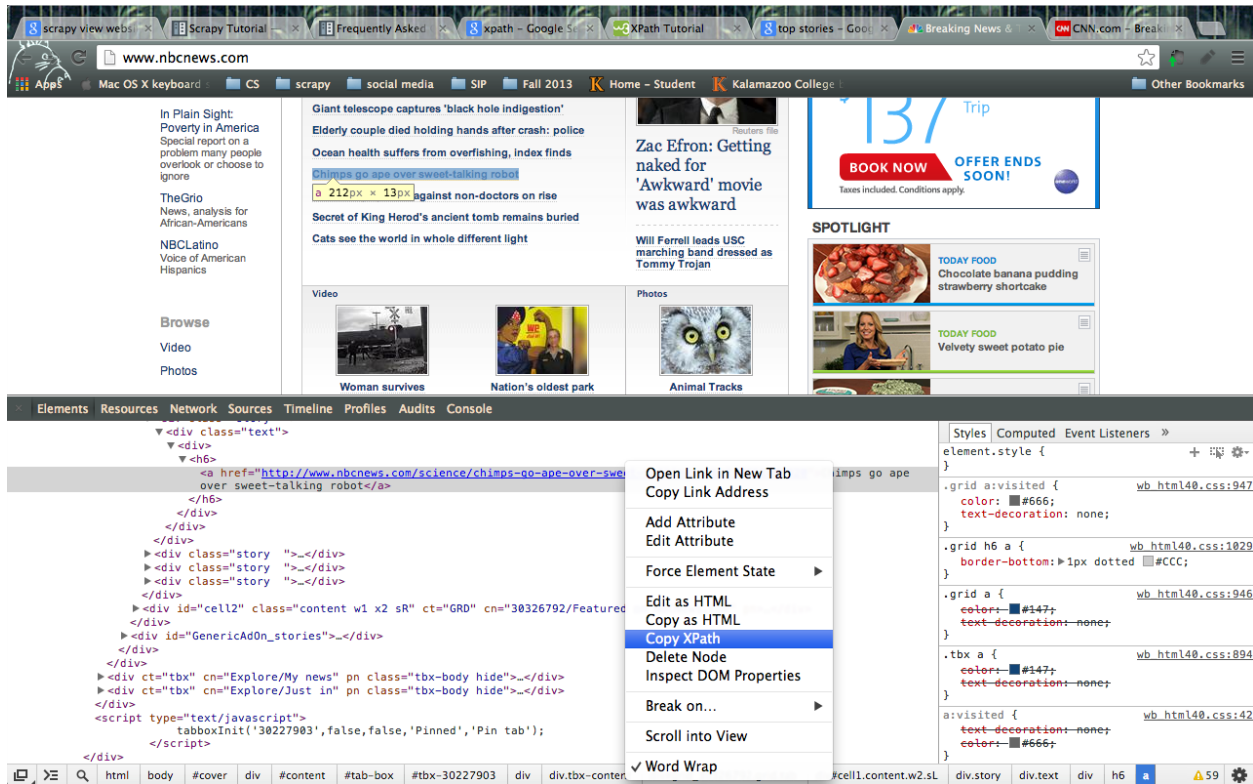


Figure 3.2: A screenshot of how to obtain an XPath from a webpage, and the XPath of the link is `//*[@id="cell1"]//div[8]/div/div/h6/a`.

```

class Spider(BaseSpider):
    def parse(self, response):
        ##some code to extract data from this current page
        yield reslts_from_this_page
        if next_page:
            yield Request(url_of_the_next_page,
                           callback=self.parse)

```

In the code above, *BaseSpider* is one of the spider classes the user has to inherit to build a user-defined spider. In any spider class, a method like *parse(self,response)* is necessary to communicate with the computer and instruct which fields in the html to scrape from. The *parse* method yields a generator, and the code does not necessarily end there. In the event that there is a next page with the same html format, the parse method will be able to go into the new html webpage and collect the new data in the same way as it did for the previous page. Therefore, the method *parse* can also yield a *Request*. *Request* represents an HTTP Request in many cases. It sends a signal to the scraper and asks for a *Response*, which will be handed to the spider and the spider will scrape it.

3.4 Specification

After sketching and testing the basic functionalities, it is time to study what types of information the clients want, and what types of information are needed to make database management easier.

The clients usually ask for specific data that are useful to their research. For instance, an author may want to read all of his books' reviews. To collect book reviews from websites such as <http://amazon.com> and <http://goodreads.com>, we need different spiders for different websites, but potentially generate tables with identical column names so that the

author can view everything in one sheet. That data sheet may contain book name, edition number, publisher, review website, user name of the reviewer, date of review, ratings, review bodies, and other information as the column names.

On the developer's end, however, the structure will be very different. To track every piece of review, the developer may need a column with the url from which that particular review was pulled. To establish an iterator, the developer may need to track the book identification number on the website to generate the correct url to scrape from. To deal with different encodings and web structures of different websites, the developer needs to customize every spider.

A piece of pseudo-code for extracting a certain text could be:

```
item['some_field']=HtmlXPathSelector(response).select('SOME/  
XPATH/IN/HTML/text()').extract()
```

3.5 Produce a sample output

After studying what and how the spider should scrape from webpages, we are ready to produce a sample output. The sample output is necessary for three reasons.

First, formatting and encoding. It is not unusual when the ideal is far from reality. Although XPath offers a specific location of all scrapable information, it is highly possible that the formatting is not desirable. Also, most websites are written in utf-8, but Python's csv module defaults encoding to ascii, a paragraph often has a `\u` in front of it. For instance, a review "I like this!" will be translated into `\u" , , , I like this!"` through Scrapy.

Second, communication with clients and manager. Only when the sample data sheet is out, can the results of scraping be evaluated. There are limitations to scraping, so it is essential to make clients and managers aware of those. For instance, the time of producing a sample output influences the frequency of running the scraper. In the event that the target

	A	B	C	D	E	F	G	H
1	source	star	user	title	date	curpage	review	
2	Amazon	3.0 out of 5	s Lau Velazque	The Kindle Ec	5-Nov-12	4	The book itself is fab	
3	Amazon	5.0 out of 5	s Camille K.	Exactly what	5-Nov-12	4	I've been reading Sm	
4	Amazon	4.0 out of 5	s Akiko	A Good Cook	6-Nov-12	4	In the last 30 days, I	
5	Amazon	5.0 out of 5	s Suzanne Trar	Smitten Kitch	6-Nov-12	4	My 4 year old daught	
6	Amazon	5.0 out of 5	s Rebekah	love the blog	6-Nov-12	4	I adore Deb's blog, I\	
7	Amazon	5.0 out of 5	s KRRRI "akat04	Yum, yum, yu	6-Nov-12	4	I have followed the S	
8	Amazon	5.0 out of 5	s Samantha Sa	Everything I e	6-Nov-12	4	Couldn\'t wait for thi	
9	Amazon	2.0 out of 5	s Love cookbo	Not enough e	6-Nov-12	4	As much as I love the	
10	Amazon	5.0 out of 5	s SamanthaJes	A deep, forev	7-Nov-12	4	I love this cookbook t	

Figure 3.4: What extracted data may look like in csv format.

scrapes the 50 pages on day 1, it store the last page as 50, so that when there are 60 pages on day 2, the scraper can start scraping from page 50 to 60, without re-scraping the first 49 pages.

For another instance, one can make the scraper more interactive with the user. One can use Pipeline to save the output data into a user specified directory. For some scraping job like product reviews, it is also possible to ask for an identifier of the product, and the scraping process will start from an auto-generated url.

Chapter 4

Future Development of the Project

4.1 Natural Language Toolkit

After harvesting data from online sources, there are various ways to enjoy the fruit. This chapter will illustrate some of possible ways to use the data.

Natural Language ToolKit, or NLTK, has many functions, and one of them is tokenization. Tokenization is the process of breaking up a paragraph into sentences, a sentence into phrases, a phrase to words, or even a word to its components. These subsets are called tokens, and more study can be done through analyzing these tokens.

For instance, if a product manager would like to monitor the sentiment of customers on the product, he/she could use web-scraping to collect review texts from online rating websites. Although those texts are available, the feedback from the reviews are not quantitative. With NLTK, however, one can tokenize each review into sentences, and sentences into phrases. For example, one can assign a score of “1” if the review says “I hate this product!”, and a score of “10” if the review reads “I love this product!”. There are databases such as WordNet available for download as a package for NLTK, and there are algorithms for assigning scores to words.

However, there are few challenges in the process of detecting sentiment.

First, the ambiguity of language. Same word may have different meanings in different context. For instance, in the review of a product, “I am sick of it” and “This product is sick!” have opposite sentiments, but they both use the word “sick” to express the sentiment. As the internet language becomes increasingly versatile, with slangs, emoticons, puns, and other forms of language in the scraped text, the difficulty of guessing the right sentiment of an expression increases drastically.

Second, the complexity of language. Short statements such as “This product is great!” and “I hate it!” are generally easier for sentiment detection. As the statement becomes longer into a complex sentence, a paragraph or even multiple paragraphs, the evaluation of a certain review becomes more complicated, too. For instance, for a review such as “Though this product may not be the worst in the world, please choose other products if possible.”, it seems obvious to human that the sentence contains negative sentiment against the product. To the computer, however, when it sees “not” followed by “worst”, it may assign a very positive sentiment to this review.

4.2 Data Mining

Web-scraping has the potential of collecting data from multiple facets. When a review website displays reviews on its webpage, information such as date, ratings, review texts and other relevant details will be available as well.

While it seems feasible to operate data mining under many collectable variables, the actual practice will not be easy. There are two main limitations.

First, limitation of parsing. The default format of scraped data is string, but more often than not, numbers are preferred in data mining. The process of parsing numbers from their string representations may incur errors.

Second, limitation of accessible data. Websites are not likely to display some data, such as demographic data, to the public. To investigate possible correlations, insufficient data will be available for scraping.

4.3 Graphic User Interface

For future development of web-scraping, it is possible to build an user interface from the code.

To scrape information related to a particular item, the user can choose a website to scrape from, and input necessary fields such as the unique identifier of that item in the website. Then, the user can choose from a checklist of what areas to scrape from, and the number of results shown. Then, after clicking a button, the Python script at the back end will use Scrapy to retrieve the input information, create starting and ending urls, and scrape selected results from relevant pages.

However, there are few limitations of establishing an interface.

First, validity of user input. As the back end Python script retrieves information from the user interface, invalid input may cause the program to stop functioning.

Second, constantly changing web interfaces. As web-scraping depends heavily on web structures, in the event of a change in web interfaces, the spider will probably fail to crawl information. Therefore, both the programmer and the user need to update frequently. The stability of the software will be very fragile.

Appendices

Glossary

AJAX Asynchronous JavaScript and XML. To quote from w3schools, “AJAX allows web pages to be updated asynchronously by exchanging small amounts of data with the server behind the scenes. This means that it is possible to update parts of a web page, without reloading the whole page.” For more information, visit http://www.w3schools.com/ajax/ajax_intro.asp 8

API Application Programming Interface. It defines and explains how users can use components of a program. A web’s API usually guides users into permitted ways of extracting information from its database. 7, 12

Business Intelligence technologies to congregate data for business analysis. 11

crawl the process of downloading webpages and copy certain information from them. 6

csv Comma-Separated Values, a plain text form of tabulated data. Entries in the same row are separated by commas,(sometimes pipelines or other characters) and columns are separated by newlines. 5, 6

Data Mining The practice of searching through large amounts of computerized data to find useful patterns or trends. (Definition from Merriam-Webster) 7, 11

encoding and decoding The conversion between elements of a certain language and numbers and texts. For instance, the character “A” is encoded as “0X41” in UTF-8 hex format, and “01110001” in UTF-8 binary format decodes to the character “q”. 9

harvest Data Harvesting, the automated process of gathering and organizing data. 5, 7, 12

HTTP Request A message sent from the client to the server through actions such as clicking a button on a webpage, asking for information. The server will examine the message, and send a response message back to the client. 19

JavaScript A computer programming language that allows users to interact with contents on webpages. 8, 17

json JavaScript Object Notation, a human-readable format for storing data. It consists mostly of name/value pairs, so different parts of a data table will be easily recognized by their names. 5, 12

NLP Natural Language Processing, the study of using computers to interpret natural languages through linguistic analysis. 11

NLTK Natural Language ToolKit, a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to over 50 corpora and lexical resources such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning.(Extracted from <http://nltk.org/>) 23

XMLHttpRequest A JavaScript object. It provides an easy way to retrieve data from a URL without having to do a full page refresh, allowing the response to be loaded within the script.

(From <https://developer.mozilla.org/en-US/docs/Web/API/XMLHttpRequest>)

8