

Web Scraping

Introduction

- Gathering of data from the internet
- Also known as screen scraping or web harvesting (*or data mining in a trivial sense*)
- Practice of gathering data by writing an automated program that:
 - Queries a web server
 - Requests data (HTML and other files that comprise web pages)
 - Parses that data to extract meaningful information

Why Web Scrapers?

- Internet is awash in data, browsers are handy to “browse” through information but not to collect and sift through data
- Web scrapers are excellent at gathering and processing large amounts of data
- Data is often unstructured or semi-structured
- Data often isn’t made available as a formal data set
- Scraping opens up a new world of data to researchers

Web Scraping Data

- Publicly accessible data that doesn't come neatly packaged as a formal data set ...
 - online classified ads and rental housing data
 - social media and behavioral data
 - discussion forums and sentiment analysis
 - auction sites and retail price data
- Leading to extremely practical applications
 - market forecasting
 - medical diagnostics
 - natural/machine language translation
- Is Web Scraping legal?!

Legal Considerations

- Procedurally similar to browsing and indexing:
Is it publicly available information?
- Do you need to log in to access it?
- Do the terms of service explicitly forbid scraping?
- Are you using the data in a way that harms the source?
- Ethical questions
 - Can you partner with the data source's organization?
 - Proper attribution

Terminology

- **Spider** – crawls the web by following links
- **Crawler** – just another name for a spider
- **Data Scraper** – generic computer program that extracts human-readable data
- **Web Scraper** – a data scraper specifically for web pages
- **Internet bot/web robot**, or simply **bot** is a software program that runs automated tasks over the Internet

Web Scraper specifics

- Small computer program that: accesses web pages
- Finds specified data elements on the page
- Extracts them (and transforms them if necessary)
- Compiles this data into a coherent data set

Web browser Vs Scraper

- Compare Scraper's behavior to that of a web browser
- Scraper can:
 - be run iteratively over many web pages
 - access data spread across thousands or millions of pages
 - construct large, robust data sets out of otherwise messy text that would only appear in your web browser

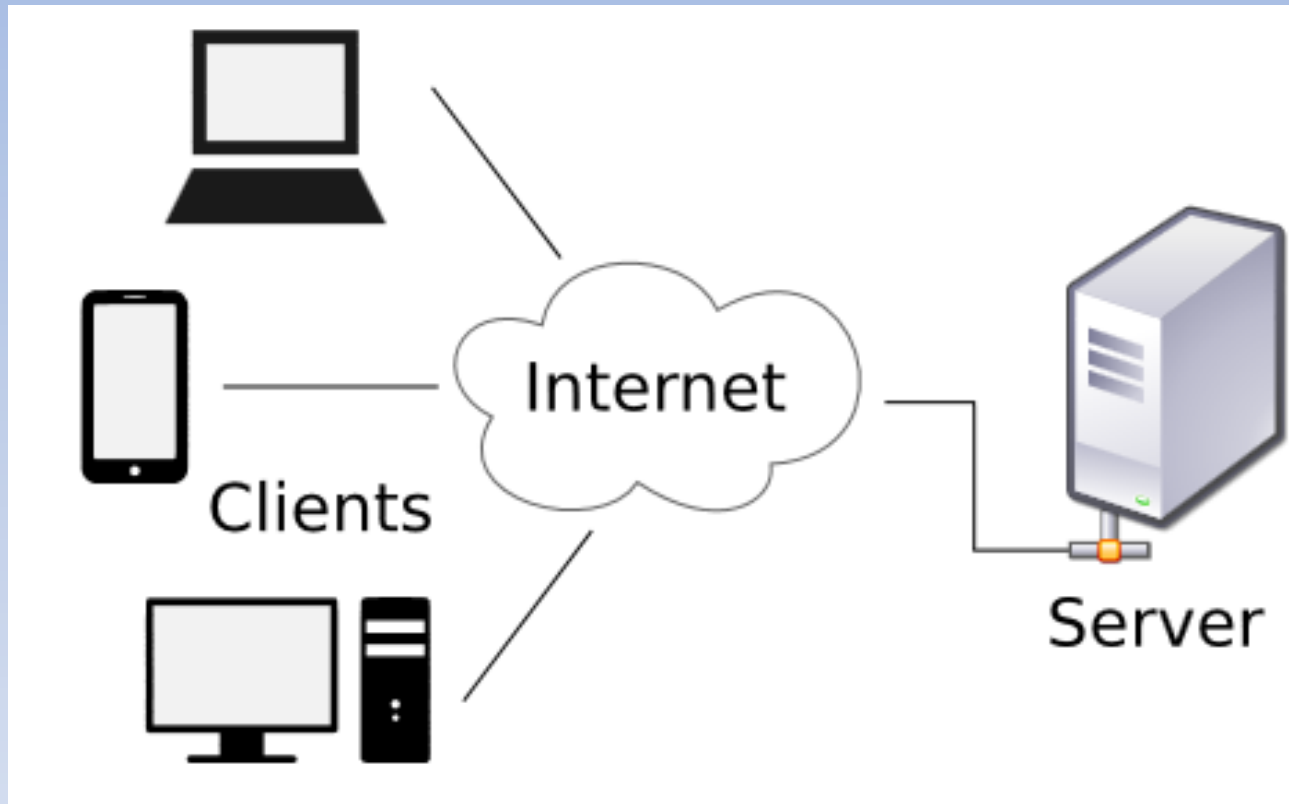
Web Scraping basics

- Traditional data transfer of unambiguous structured data
 - not very human-readable
- APIs application programming interfaces
 - provides external access to some software, data, or service
 - not all applications have APIs
 - Web scraper can be written as a custom API for some web site of interest
- Web pages
 - unstructured or semi-structured data
 - very human-readable
 - Challenge is to convert semi-structured data on web pages into a structured data set on your hard drive

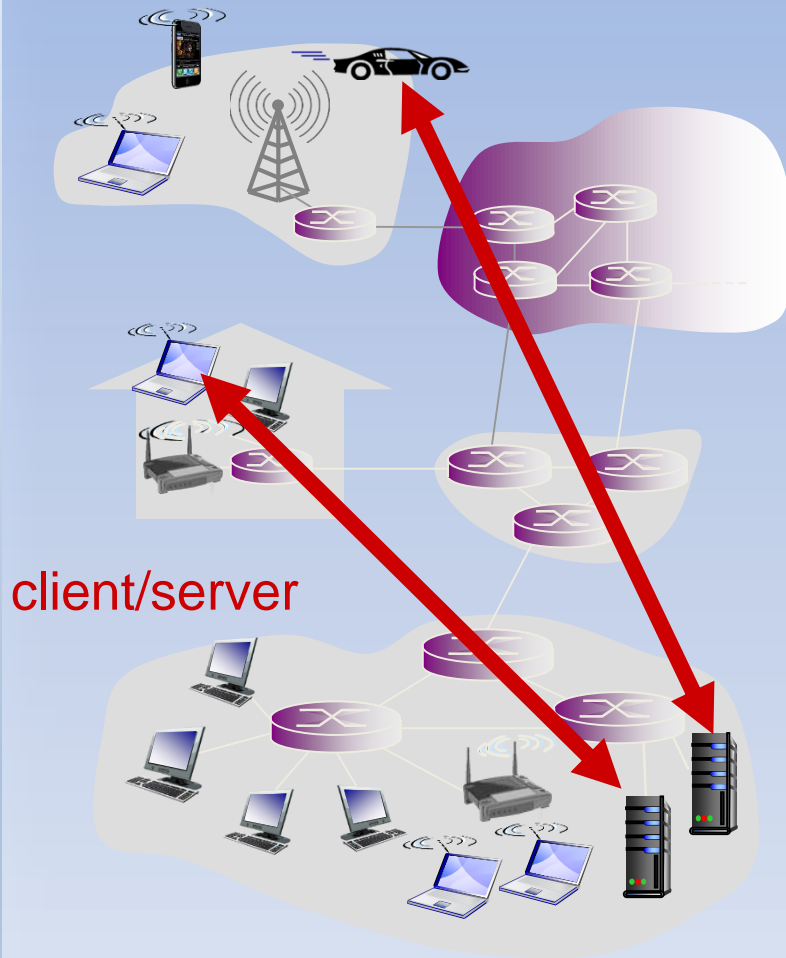
Web Scrapping basics



Client Server Model



Client-server architecture



server:

- always-on host
- permanent IP address
- data centers for scaling

clients:

- communicate with server
- may be intermittently connected
- may have dynamic IP addresses
- do not communicate directly with each other

Web and HTTP

- *web page* consists of *objects*
- object can be HTML file, JPEG image, Java applet, audio file,...
- web page consists of *base HTML-file* which includes *several referenced objects*
- each object is addressable by a *URL*, e.g.,
`www.someschool.edu/someDept/pic.gif`

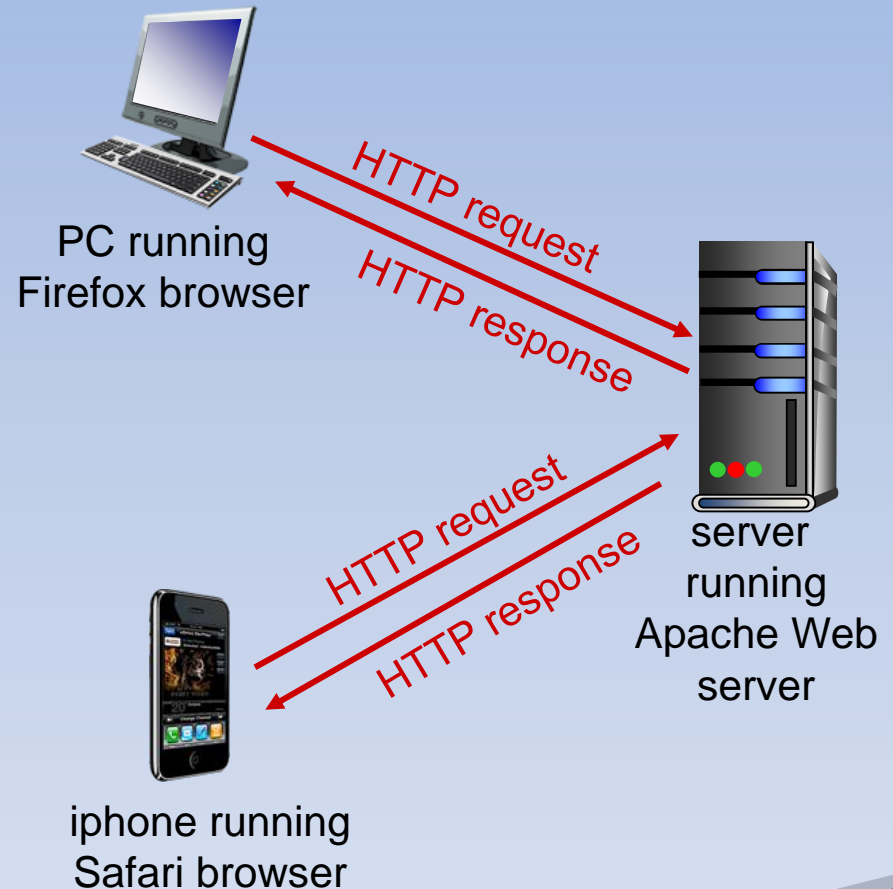
host name

path name

HTTP overview

HTTP: hypertext transfer protocol

- Web's application layer protocol
- client/server model
 - *client*: browser that requests, receives, (using HTTP protocol) and “displays” Web objects
 - *server*: Web server sends (using HTTP protocol) objects in response to requests



Semi-Structured data

> boston/camb/brook > housing > real estate - by owner [account]

search real estate - by owner

list thumb gallery map << < prev 1 to 100 of 536 next > newest \$\$\$ \$\$\$\$

- ★ Sep 26 [FORMER "NORMAN ROCKWELL GALLERY](#) - \$135000 / 3br - 3800ft² - (ARLINGTON VT) [pic](#) [map](#) [x]
- ★ Sep 26 [Looking for Cash Investors](#) - [x]
- ★ Sep 26 [MANUFACTURED HOME-DOUBLE WIDE](#) - \$72000 / 4br - 1344ft² - (HOMOSASSA) [pic](#) [map](#) [x]
- ★ Sep 26 [LAWN MOWING, PRUNING, GRASS SEEDING & MORE LANDSCAPING SERVICES](#) - [x]
- ★ Sep 26 [Lender/Investors Wanted for Real Estate Projects - 20% Return](#) - (Boston) [map](#) [x]
- ★ Sep 26 [2+ Bedroom 2 Bathroom Walk to Lake Winnepesaukee](#) - \$72000 / 2br - 1100ft² - (Laconia) [map](#) [x]
- ★ Sep 26 [MAINE HOME BUILDER](#) - (MAINE) [map](#) [x]
- ★ Sep 26 [LAND for SALE at BEACH](#) - \$35000 (Boiling Springs Lakes N.C.) [map](#) [x]
- ★ Sep 26 [condo in Florida](#) - \$49000 / 1br - (Delray Beach) [pic](#) [map](#) [x]
- ★ Sep 26 [LESS THAN 2 HRS, TO THIS GREAT HOME,WITH 3.5 ACRES NEAR WEBSTER LAKE](#) - \$124900 / 2br - 950ft² - (280 LAKE SHORE DR. FRANKLIN,N.H.)
- ★ Sep 26 [Raw Land FSBO](#) - (Essex) [map](#) [x]

- Human-readable listing
- How do we see or capture all of this rental market data?
- Right-click and view source...

Web page source code

```
1 <!DOCTYPE html>
2
3 <html class="no-js"><head>
4   <title>boston real estate - by owner - craigslist</title>
5
6   <meta name="description" content="boston real estate - by owner - craigslist">
7   <link rel="canonical" href="https://boston.craigslist.org/search/gbs/reo">
8   <link rel="alternate" type="application/rss+xml" href="https://boston.craigslist.org/search/gbs/reo?format=rss" title="RSS feed for
craigslist | boston real estate - by owner - craigslist ">
9
10  <link rel="next" href="https://boston.craigslist.org/search/gbs/reo?s=100">
11  <meta name="viewport" content="initial-scale=1.0, user-scalable=1">
12  <link type="text/css" rel="stylesheet" media="all" href="//www.craigslist.org/styles/cl.css?v=0b04254c2a971c33c5543fa6344f7dad">
13  <link type="text/css" rel="stylesheet" media="all" href="//www.craigslist.org/styles/tocs.css?v=5f8a8e4c2f016a0da7b16fc9394a62e9">
14  <link type="text/css" rel="stylesheet" media="all" href="//www.craigslist.org/styles/jquery-ui-1.9.2.custom.css?v=
86db72e37d2cf63cab3c63485c71becb">
15  <link rel="prefetch" href="//www.craigslist.org/js/postings-concat.min.js?v=2c267741ac82581ef0aff79d8144768b">
16
17  <link type="text/css" rel="stylesheet" media="all" href="//www.craigslist.org/styles/leaflet-stock.css?v=
969bf0c010aad78a472cb96ed5cd1bc">
18  <link type="text/css" rel="stylesheet" media="all" href="//www.craigslist.org/styles/MarkerCluster.css?
v=c9937ed03fbb57f185493cd8c283efeb">
19  <link type="text/css" rel="stylesheet" media="all" href="//www.craigslist.org/styles/tocsmaps.css?v=0713fba1eea156d1e29815594450b352">
20  <!--[if IE]>
21    <link type="text/css" rel="stylesheet" media="all" href="//www.craigslist.org/styles/tocsmaps-ie.css?
v=b21742b2ebe243963373e956ea115db0">
22  <![endif]>
23
24    <script type="text/javascript"><!--
25      var expiredFavIDs = [];
26  var subarea = "gbs";
```

HTML

- HTML the markup language for web pages
 - consists of HTML elements defined by tags
 - elements can have attributes and can encapsulate text
 - creates a navigable, structured document
- Element = individual component of HTML
- Elements represent semantics
- Written with a start and end tag Tags use angle brackets
- Example: `<title>this is a title</title>`
- Elements can have attributes
 - Example: `<p class="blog">this is a paragraph</p>`