

# Obtaining, Cleaning, and Exploring Corpora

Elliott Ash

August 26, 2018

# Publicly Available Corpora

- ▶ There is already a vast amount of data out there that has already been compiled (e.g. CourtListener, Twitter, New York Times, Reuters, Google, Wikipedia).
- ▶ Chris Bail curates a list of these datasets:
  - ▶ <https://docs.google.com/spreadsheets/d/1I7cvuCBQxosQK2evTcdL3qtglaEPc0WFEs6rZMx-xiE/edit>
- ▶ Some interesting corpora described in NLTK Book Chapter 2.
- ▶ Many proprietary corpora are becoming available for research:
  - ▶ Lexis
  - ▶ Web of Science

# Screen Scraping

- ▶ A screen scraper is a computer program that:
  - ▶ loads/reads in a web page
  - ▶ finds some information on it
  - ▶ grabs the information
  - ▶ stores it in a dataset
- ▶ Once upon a time you could collect virtually any piece of information from the internet by screen scraping.
  - ▶ But now web sites make it difficult with restrictive terms of use, bot-blockers, javascript, etc.
  - ▶ Still, a little creativity goes a long way.

# What a web site looks like to us



WIKIPEDIA  
The Free Encyclopedia

[Main page](#)  
[Contents](#)  
[Featured content](#)  
[Current events](#)  
[Random article](#)  
[Donate to Wikipedia](#)  
[Wikipedia store](#)

Interaction

[Help](#)  
[About Wikipedia](#)  
[Community portal](#)  
[Recent changes](#)  
[Contact page](#)

Tools

[What links here](#)

Create account Log in

Article [Talk](#)

Read [Edit](#) [View history](#)

## World Health Organization ranking of health systems in 2000

From Wikipedia, the free encyclopedia

The **World Health Organization (WHO)** **ranked the health systems** of its 191 member states in its [World Health Report](#)<sup>[1]</sup> 2000. It provided a framework and measurement approach to examine and compare aspects of [health systems](#) around the world.<sup>[2]</sup> It developed a series of performance indicators to assess the overall level and distribution of [health](#) in the populations, and the responsiveness and financing of [health care](#) services. It was the organization's first ever analysis of the world's health systems.<sup>[3]</sup>

**Contents** [\[hide\]](#)

- [1 Ranking](#)
- [2 Methodology](#)
- [3 Criticism](#)
- [4 See also](#)
- [5 References](#)

# What a web site looks like to a computer

```
1 <!DOCTYPE html>
2 <html lang="en" dir="ltr" class="client-nojs">
3 <head>
4 <meta charset="UTF-8" />
5 <title>World Health Organization ranking of health systems in 2000 - Wikipedia, the free encyclopedia</title>
6 <meta name="generator" content="MediaWiki 1.26wmf10" />
7 <link rel="alternate" href="android-
  app://org.wikipedia/http/en.m.wikipedia.org/wiki/World_Health_Organization_ranking_of_health_systems_in_2000"
  />
8 <link rel="alternate" type="application/x-wiki" title="Edit this page" href="/w/index.php?
  title=World_Health_Organization_ranking_of_health_systems_in_2000&action=edit" />
9 <link rel="edit" title="Edit this page" href="/w/index.php?
  title=World_Health_Organization_ranking_of_health_systems_in_2000&action=edit" />
10 <link rel="apple-touch-icon" href="/static/apple-touch/wikipedia.png" />
11 <link rel="shortcut icon" href="/static/favicon/wikipedia.ico" />
12 <link rel="search" type="application/opensearchdescription+xml" href="/w/opensearch_desc.php" title="Wikipedia
  (en)" />
13 <link rel="EditURI" type="application/rsd+xml" href="//en.wikipedia.org/w/api.php?action=rsd" />
14 <link rel="alternate" hreflang="x-default"
  href="/wiki/World_Health_Organization_ranking_of_health_systems_in_2000" />
15 <link rel="copyright" href="//creativecommons.org/licenses/by-sa/3.0/" />
16 <link rel="alternate" type="application/atom+xml" title="Wikipedia Atom feed" href="/w/index.php?
  title=Special:RecentChanges&feed=atom" />
17 <link rel="canonical"
  href="https://en.wikipedia.org/wiki/World_Health_Organization_ranking_of_health_systems_in_2000" />
18 <link rel="stylesheet" href="//en.wikipedia.org/w/load.php?
  debug=false&lang=en&modules=ext.uls.nojs%7Cext.visualEditor.viewPageTarget.noscript%7Cext.wikihihero%7C
  mediawiki.legacy.commonPrint%2Cshared%7Cmediawiki.sectionAnchor%7Cmediawiki.skinning.interface%7Cmediawiki.ui.
  button%7Cskins.vector.styles%7Cwikibase.client.init&only=styles&skin=vector&*"/>
19 <meta name="ResourceLoaderDynamicStyles" content="" />
20 <link rel="stylesheet" href="//en.wikipedia.org/w/load.php?
  debug=false&lang=en&modules=sites&only=styles&skin=vector&skin=vector&*"/>
https://en.wikipedia.org/w/index.php?title=World_Health_Organization&a:lang(mzn),a:lang(ps),a:lang(ur){text-decoration:none}
```

# Screen Scraping in Python

```
# package to access web pages  
from urllib.request import urlopen  
url = 'https://goo.gl/VRF8Xs' # shortened URL  
page = urlopen(url) # open page  
  
html = page.read() # read page as string  
print(html[:400]) # print first 400 characters  
print(html[-400:]) # print last 400 characters  
print(len(html)) # print length of string
```

# Browser Automation

- ▶ Many web sites are designed to be difficult to scrape.
- ▶ Python has many solutions for simulating a human browser:
  - ▶ robobrowser
  - ▶ selenium (chromedriver, phantomjs)
- ▶ Other solutions if all else fails:
  - ▶ DownThemAll! plug-in for Firefox
  - ▶ Hire mechanical turkers to manually download data.

# API's

- ▶ API = Application Programming Interface
  - ▶ These are developer-oriented tools that provide access to cleaner data.
- ▶ Chris Bail's list of API's that could be interesting for research:
  - ▶ <https://docs.google.com/spreadsheets/d/1ZEr3okdlb0zctmX0MZKo-gZKPsq5WGn1nJ0xPV7a1-Q/edit>
- ▶ The example data set was obtained from the CourtListener API ([courtlistener.com/api](http://courtlistener.com/api)).



## Other Languages

- ▶ All of the tools that we discuss in this class are available in many languages.
- ▶ spaCy has full functionality in English, German, Spanish, Portuguese, French, Italian, and Dutch.
  - ▶ beta functionality in dozens of other languages including Chinese and Arabic
  - ▶ See <https://spacy.io/usage/models>.

The machine learning models are language-independent.

# Translation

- ▶ Can also translate languages before running an analysis

```
from googletrans import Translator
translator = Translator()
lang = translator.detect(korean_string).lang
eng = translator.translate(korean_string ,
                           src=lang ,
                           dest='en')

eng.text
```

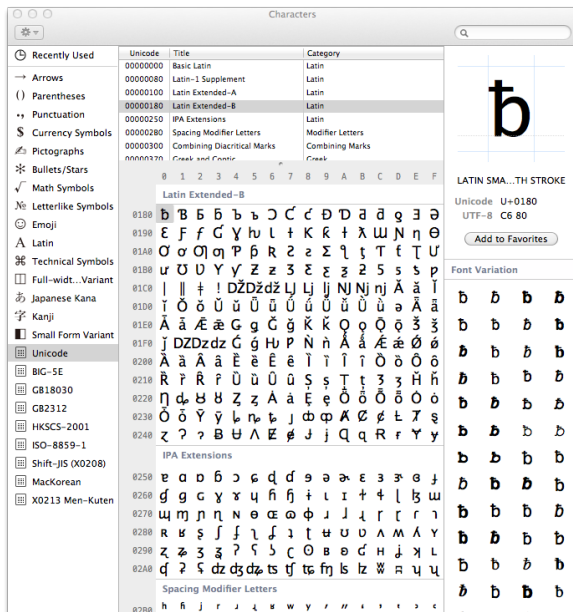
# HTML Parsing

```
from bs4 import BeautifulSoup # HTML parser
soup = BeautifulSoup(html, 'lxml') # parse page
print(soup.title)

# extract text
text = soup.get_text() # remove HTML markup
lines = text.splitlines() # split by line
print(len(lines)) # print number of lines

# drop empty lines
lines = [line for line in lines if line != '']
print(len(lines))
print(lines[:20]) # print first 20 lines
```

# Character Encoding



The screenshot shows the macOS 'Characters' application. The 'Unicode' tab is selected, displaying a grid of characters. The 'Recently Used' sidebar on the left lists various symbol categories. The main grid shows characters from different encodings, with 'Latin Extended-B' currently selected. The character 'b' is highlighted, and its details are shown on the right.

Unicode	Title	Category
00000000	Basic Latin	Latin
00000080	Latin-1 Supplement	Latin
00000100	Latin Extended-A	Latin
00000180	Latin Extended-B	Latin
00000250	IPA Extensions	Latin
00000280	Spacing Modifier Letters	Modifier Letters
00000300	Combining Diacritical Marks	Combining Marks
00000320	Greek and Coptic	Greek

Character details for 'b':

- Font Variation: LATIN SMALL LETTER B
- Unicode: U+0180
- UTF-8: C6 80
- Add to Favorites

## Removing unicode characters

```
# package for removing unicode  
from unicodecode import unicodecode  
fixed = unicodecode('Visualizations\xa0')  
print(fixed) # print cleaned string
```

# Corpus cleaning

- ▶ What we've already done:
  - ▶ removed HTML markup, extra white space, and unicode
- ▶ But HTML markup is often valuable:
  - ▶ HTML markup for section header names.
  - ▶ Legal database web sites often have HTML tags for citations to other cases.
- ▶ Other cleaning steps:
  - ▶ page numbers
  - ▶ hyphenations at line breaks
  - ▶ table of contents, indexes, etc.
- ▶ These are all corpus-specific, so inspect ahead of time.

# Regular Expressions

- ▶ Regular Expressions, implemented in the Python package **re**, provide a powerful string matching tool.
  - ▶ A systematic string matching protocol – can match arbitrary string patterns
  - ▶ e.g., use `utilit*` to match `utility`, `utilities`, `utilitarian`, ...
  - ▶ Important for identifying speaker names (in political documents) section headers (in statutes), citations (in judicial opinions), etc.
- ▶ Also quite tedious, so we will not cover it here.
  - ▶ See NLTK book Chapter 3.4-3.5 for an introduction.

# OCR (Optical Character Recognition)

- ▶ Your data might be in PDF's or images. Needs to be converted to text
- ▶ The best solution (that I know of) is ABBYY FineReader, which is expensive but might be available at your university library.
- ▶ My colleague Joe Sutherland at Columbia has a nice open-source package for OCR:
  - ▶ <https://github.com/jlsutherland/doc2text>



# Spelling corrections

- ▶ Should you run a spell checker on your corpus?
- ▶ The short answer is no:
  - ▶ Most corpora have important specialized vocabulary that would be flagged by standard spell-checkers.
  - ▶ They are also very slow to run on large corpora.
  - ▶ In most empirical contexts, it's safe to assume that spelling errors (especially OCR errors) are uncorrelated with treatment assignment.
- ▶ Better solutions:
  - ▶ drop short (one or two letters) and long words (over 12 letters).
  - ▶ get document frequencies for each word and filter out rare words
- ▶ But:
  - ▶ There are cases where spelling errors could be correlated with treatment (for example, increasing legislator salaries might change both policy priorities and spelling error rates)
  - ▶ Check out the **enchant** module in Python.

## Collect Key Metrics

```
df1 = df1[['state', 'snippet']]
# Number of documents
len(df1['snippet'])
# Number of label categories (e.g. states)
df1['state'].describe()
# Number of samples per class
counts_per_class = df1.groupby('state').count()
counts_per_class.head()
# Words per sample
def get_words_per_sample(txt):
    return len(txt.split())
df1['num_words'] = df1['snippet'].apply(get_words_per_sample)
df1['num_words'].describe()
# Frequency distribution over words
from collections import Counter
freqs = Counter()
for i, row in df1.iterrows():
    freqs.update(row['snippet'].lower().split())
freqs.most_common()[:20]
# (Number of samples) / number of words per sample
len(df1['snippet']) / df1['num_words'].mean()
```