

# Wikipedia Degrees of Separation

*Andre Duarte, Nick Levitt, Shivakanth Thudi, Juan Pablo Oberhauser*

*October 10, 2016*

## Contents

<b>Project Description</b>	<b>2</b>
<b>Data Acquisition</b>	<b>4</b>
Overview of Available Options . . . . .	4
Drawbacks and Advantages of Acquisition Options . . . . .	4
<b>CURL Solution for ‘Build a Path’</b>	<b>5</b>
Solution for List of Links . . . . .	6
<b>Solution for Degrees of Separation and ‘Find a Path’</b>	<b>7</b>
Journey through shortest path . . . . .	9
<b>Other tools used</b>	<b>10</b>
Encoding Issues . . . . .	10
Flask, Selenium, and BeautifulSoup . . . . .	10
<b>Bibliography</b>	<b>10</b>

## Project Description

The goal of this project is to create and measure a relationship in terms of degrees of separation between certain terms in Wikipedia. The degrees of separation between two Wikipedia search terms is defined by the number of clicks necessary to go from one article to the other, only using links within each subsequent Wikipedia page. This calculation is a proxy for “how closely related” those two search terms are.

Two methods of traversal can be envisioned. The first is to start at a specific Wikipedia page, and iteratively navigate Wikipedia without knowing the final landing page. The user starts at a page and clicks links until he decides to stop. The second is to find a path (hopefully the shortest one) between two terms. The number of pages in that path (minus 1) is the degrees of separation between the two terms.

In this project, we implemented both approaches to this problem. Each part of the project presents its own algorithmic, data storage, and data retrieval challenges that will be further discussed. The first option is called “**Build a Path**”, while the second is called “**Find a Path**”.

Ultimately, the user has the choice of which traversal method he wishes to explore, and so he is given a choice right from the landing page of the server, as shown in Figure 1.

### Select a traversal method

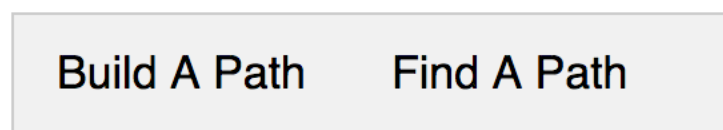


Figure 1: Landing page with 2 traversal method options

Let’s suppose a user chooses to explore Wikipedia from a starting page, with no final page in mind. He selects to start at **Italian cuisine**, which leads him to **Italians**, then **Canda**, and finally **Latin**. Figure 2 shows a glimpse of what the user sees and how the pages he has visited are interconnected.

Now, let’s imagine the user wants to find the shortest path (within a subsection of Wikipedia, more on this later) between the terms **Ageing** and **Abiogenesis**. By selecting the **Find a Path** option with these terms, he is shown the result in Figure 3.

Click here to go back to traversal selection



Italian cuisine --> Italians --> Canada --> Latin

**Latin: 731 links**

[Latin \(disambiguation\)](#)

[Romance languages](#)

[Romansh language](#)

[Romanesco dialect](#)

[Romanian language](#)

[Romani language](#)

[Colosseum](#)

[Latium](#)

[Roman Kingdom](#)

[Roman Republic](#)

[Roman Empire](#)

[Early modern Europe](#)

Figure 2: Build Path View

The shortest path from Ageing to Abiogenesis is:

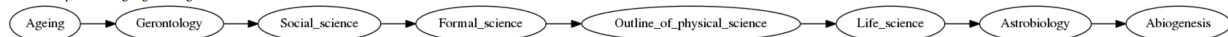


Figure 3: Shortest Path View

# Data Acquisition

## Overview of Available Options

For this project we require Wikipedia data. More specifically, we require all of the wikipedia search terms or article titles, and a list of all the links contained in each article. The wikipedia English language dump is 15 GB compressed, so we have to be creative in our solution to acquire either a subset of articles or a different methodology to get this data. We propose several solutions, their tests, and drawbacks.

1. The English language static **html dump**: “A copy of all pages from all Wikipedia wikis, in HTML form.” (14.3 GB as a .7z file)
2. **Pywikibot**: Python library that provides functionality to automate work on Wikipedia sites.
3. **SQL interconnectedness files**, provided by Wikipedia. Files contain metadata on which articles contain which links.
4. **CURL the urls** for the article that is chosen on the fly, get the list of links from a parsed HTML, and provide a choice of new articles to explore and build a navigation path.
5. Build a fully connected tree of from a **subset of articles**, limiting search options to terms included in the subset, and navigate the tree of parents (terms) and children (links in the term article) to find the shortest path.

## Drawbacks and Advantages of Acquisition Options

- The English language static **html dump**: “A copy of all pages from all Wikipedia wikis, in HTML form.” (14.3 GB as a .7z file)

This option had to be quickly discarded because even though the dump was successfully stored in our local machines, it would have been impossible to host on pythonanywhere.com. pythonanywhere.com limits the amount of storage on disk that is available to the user as well as the ‘CPU seconds’. ‘CPU seconds’ would not necessarily be the limiting factor for our project since these are only used up when the CPU is actually busy. Our limiting factor was more on the storage limitations.

- **Pywikibot**: Python library that provides functionality to automate work on Wikipedia sites.

Using a python library that provides functionality to work with wikipedia pages seemed very reasonable since it seems like many people around the web are using it to update article content and give maintenance to content.

Pywikibot allows about 5,000 article downloads without being stopped. Even though this seems like a viable solution, we need to make sure that the subset of articles is very interconnected so that there can be a shortest path solution between terms. Making sure that the subset was that interconnected proved to be a difficult task and this option had to be foregone.

- **SQL interconnectedness files**, provided by Wikipedia. Files contain metadata on which articles contain which links.

This idea proved doable with a subset of data from the Philippines. We were able to download the files and search through the content in the file. We used a shortest path algorithm and were able to always get the shortest path between two terms. When trying to switch to the full English interconnectedness files, storage again became a limiting factor. The English language full dump is 40GB, so this again proved impossible.

- **CURL the urls** for the article that is chosen on the fly, get the list of links from a parsed HTML, and provide a choice of new articles to explore and build a navigation path.

This was the viable solution for the data part of the ‘Build a Path’ part of the project.

- Build a fully connected tree of from a **subset of articles**, limiting search options to terms included in the subset, and navigate the tree of parents(terms) and children(links in the term article) to find the shortest path.

The two methods that are used in the final implementation of this project are **CURLing articles on the fly** for the ‘Build a Path’ part, and getting a subset of 400 fully-connected articles to build a traversable tree for the ‘Find a Path’ part of the project.

## CURL Solution for ‘Build a Path’

For the ‘Build a Path’ section, we use urllib2 functions, specifically urlopen(). This allows us to open a **network object** for reading. To get the URL, we get the user-input term and make a URL to send a request to. To ensure functionality, we use error handling in the URLs in the following manner.

```

def getLinksFromSearchString(searchString):
    searchURL = template_wikiURL % searchString
    URL = mainWikiURL % searchURL
    try:
        return getLinksFromURL(URL)
    except:
        return ["%s is not a valid search term" % searchString]

```

## Solution for List of Links

Given the difficulty of dealing with a data dump and the difficulty of using the python bot to extract information, we will be implementing a script that does a request to a given URL and brings back the html, to be parsed using BeautifulSoup.

After the retrieval of the text, we look at all of the pieces of html containing a “class: mw-content-ltr”

For example:

```

<div class="mw-content-ltr" dir="ltr" id="mw-content-text" lang="en">
<div class="hatnote" role="note">
This article is about artificial satellites. For natural satellites, also known as moons,
see <a href="/wiki/Natural_satellite" title="Natural satellite"> Natural satellite</a>

```

Once we obtain a link object, we look for the ‘a’ anchor to get the actual part of the text which contains the ‘href’ , we store that link and its name in a python list to be utilized when the html is rendered. At this stage of the search for a list of links in the article body, it is necessary to check for a valid link. To do so, we use string properties to check whether the URL string starts with ‘wiki/’ and check for invalid characters. This way, for each page that is scraped, we store all the links to other Wikipedia pages in a list.

## Storing linked elements

As a first solution to store the connectedness of the articles, we thought about using a linked list approach. This would allow us to begin with a node (which is the first search term as input from the user) and a pointer, which will reference the next node on the list. The next node on the list will then contain the name of the article that was clicked from that first list, and so on.

Although this method seems intuitive, it would not work for our purposes. In fact, we wish to store two pieces of information: the sequential list of articles visited, as well as how they connect with each other. If the user went back to a page that he previously visited, then the order of the visits would get lost.

We therefore opted to keep two structures to store the information. First, a list of visited pages, where each subsequent page is appended to the list. Second, a dictionary of connected pages, where each key is associated with the article that it links to.

Another advantage of using a dictionary in this way is that it can easily be transformed into a directed graph text notation that can be passed to the command line's **dot** function, which produces a graph image (that we display to the user).

## Solution for Degrees of Separation and ‘Find a Path’

As opposed to the ‘Build a Path’ section of the project, the ‘Find a Path’ needs to find the shortest path between two given terms that are user-input. This provides a challenging data acquisition project since to find a path from Article A to Article B, we need to ensure there is in fact a way to get there using only links in each article.

The proposed solution for this problem is the following:

1. Get a subset of 400 articles (Parent)
2. For each article, get 5 links contained in the article (Children), making sure that each Child links back to the Parent.
3. Create a traversable tree with 400 Parent nodes, each of which has 5 children that link back to the parent.

The limit of 400 articles makes sure that we have a very interconnected set of articles. Unlike the ‘Build a Path’ part which gets data on the fly, the ‘Find a Path’ section has all of the data pre-stored on disk and a tree of nodes that will traverse and find the nearest path between two terms each time the user inputs two terms. The tree is stored in a txt file.

## Data types for storing hierarchical data

To be able to find relationships between the terms that the user chooses (i.e find the shortest path between two terms), we need a very specific data structure. To iterate a tree-like data structure, python provides list functionality.

A simple example is:

```
def walk(node):  
    """ iterate tree in pre-order depth-first search order """  
    yield node  
    for child in node.children:  
        for n in walk(child):  
            yield n
```

Figure 4: Tree-like iteration with lists



## Journey through shortest path

As a visual bonus for the user, we have included a visual journey through the shortest path between two articles. Harnessing Selenium, we can create a visual tour that mimics that shortest path and actually opens a browser window showing each click that it would take to get from Article A to Article B.

The screenshot shows a Wikipedia page for 'Taxonomy (biology)'. The 'Domains' section is highlighted, discussing the three-domain system and the work of Thomas Cavalier-Smith. To the right, a diagram illustrates the taxonomic hierarchy from Kingdom to Species. Below the text is a table comparing different taxonomic systems.

Linnaeus 1735 <sup>[32]</sup>	Haeckel 1866 <sup>[33]</sup>	Chatton 1925 <sup>[34]</sup>	Copeland 1938 <sup>[35]</sup>	Whittaker 1969 <sup>[36]</sup>	Woese et al. 1990 <sup>[37]</sup>	Cavalier-Smith 1998 <sup>[30]</sup>
2 kingdoms	3 kingdoms	2 empires	4 kingdoms	5 kingdoms	3 domains	6 kingdoms
(not treated)	Protista	Prokaryota	Monera	Monera	Bacteria Archaea	Bacteria
			Protoctista	Protista		Protozoa Chromista
Vegetabilia	Plantae	Eukaryota	Plantae	Plantae	Eucarya	Plantae
Animalia	Animalia		Animalia	Animalia		Fungi Animalia

Main article: [Kingdom \(biology\) § Summary](#)

**Application** [edit]

Biological taxonomy is a sub-discipline of [biology](#), and is generally practised by biologists known as "taxonomists", though enthusiastic [naturalists](#) are also frequently involved in the publication of new taxa. The work carried out by taxonomists is crucial for the understanding of biology in general. Two fields of applied biology in which taxonomic work is of fundamental importance are the

Figure 5: Journey Step Example

In the example above, the journey has taken us to the 'Taxonomy\_biology' page and is now highlighting what the next step will be, which in this case is 'Domains'. This specific example is a step in the journey to get to the 'Number 1' Wikipedia page from the 'Almhult Municipality' page.

Since the final project is being hosted on pythonanywhere.com, we need error handling around this 'Take the Journey' part because it would not be possible to open a driver with Selenium on this platform for each user. For this reason, we disable the take the journey part, and instead print a visual representation of what the shortest path would look like.

## Other tools used

### Encoding Issues

Given that some of the titles in the wikipedia articles contain non-ASCII characters, we provide an automated solution where each time article information is retrieved from the web, we convert everything to ASCII characters.

### Flask, Selenium, and BeautifulSoup

To get the list of links, we use **BeautifulSoup** to parse the html returned from our request to Wikipedia.

To provide the “journey” functionality, **Selenium** python bindings were utilized. These bindings provide API access to web-drivers or web browsers.

**Flask** is used to provide all of the functional bindings to turn functions into html. What Flask does is map a url to a python function. This allows us to set up event methods that are mapped to a particular URL.

## Bibliography

- Python Tree Data Sctructure. (n.d.). Retrieved from <http://stackoverflow.com/questions/3009935/looking-for-a-good-python-tree-data-structure>
- Open Arbitrary Resources by URL. (n.d.). Retrieved from <https://docs.python.org/2/library/urllib.html>
- What are “CPU seconds”?. (n.d.). Retrieved from <https://help.pythonanywhere.com/pages/WhatAreCPUSeconds>
- Data Structures. (n.d.). Retrieved from <https://docs.python.org/3/tutorial/datastructures.html>