

**Digging Deeper,
Reaching Further**

Module 3: Working with Textual Data

In this module we'll...

- Think about what happens when text is data
 - *Understand best practice in the field*
- Consider common steps to cleaning and preparing text data
 - *Make recommendations to researchers*
- Practice data cleaning on the speeches we scraped
 - *Gain experience with Python for working with data*
- Learn how Sam prepared his *Creativity Corpus* for analysis
 - *See how one scholar data prepared data*



Where we'll end up

```
1 <!DOCTYPE html>
2 <html class="client-nojs" lang="en" dir="ltr">
3 <head>
4 <meta charset="UTF-8"/>
5 <title>George Washington's Fourth State of the Union Address - Wikisource, the free online library</title>
6 <script>document.documentElement.className = document.documentElement.className.replace( /(\s)client-nojs(\s$)/, "$1client-js$2" );</script>
7 <script>(window.RLQ=window.RLQ||[]).push(function(){mw.config.set({"wgCanonicalNamespace":"","wgCanonicalSpecialPageName":false,"wgNamespaceNumber":0,"wgPageName":"George_Washington's_Fourth_Sta
8 mw.user.tokens.set({"editToken":"+\\","patrolToken":"+\\","watchToken":"+\\","csrfToken":"+\\"});/*@nomin*/;
9
10 });mw.loader.load(["site","mediawiki.page.startup","mediawiki.user","mediawiki.hidpi","mediawiki.page.ready","mediawiki.searchSuggest","ext.gadget.charinsert","ext.gadget.Easy_LST","ext.gadget.d
11 <link rel="stylesheet" href="/w/load.php?debug=false&lang=en&modules=ext.uls.interlanguage&7Cext.visualEditor.desktopArticleTarget.noscript&7Cext.wikimediaBadges&7Cmediawiki.legacy.commo
12 <script async="" src="/w/load.php?debug=false&lang=en&modules=startup&only=scripts&skin=vector"></script>
13 <meta name="ResourceLoaderDynamicStyles" content="" />
14 <link rel="stylesheet" href="/w/load.php?debug=false&lang=en&modules=ext.gadget.Site&only=styles&skin=vector"/>
15 <link rel="stylesheet" href="/w/load.php?debug=false&lang=en&modules=site.styles&only=styles&skin=vector"/>
16 <meta name="generator" content="MediaWiki 1.30.0-wmf.7"/>
17 <meta name="referrer" content="origin-when-cross-origin"/>
18 <link rel="alternate" type="application/x-wiki" title="Edit" href="/w/index.php?title=George_Washington%27s_Fourth_State_of_the_Union_Address&action=edit"/>
19 <link rel="edit" title="Edit" href="/w/index.php?title=George_Washington%27s_Fourth_State_of_the_Union_Address&action=edit"/>
20 <link rel="shortcut icon" href="/static/favicon/wikisource.ico"/>
21 <link rel="search" type="application/opensearchdescription+xml" href="/w/opensearch_desc.php" title="
22 <link rel="EditURI" type="application/rsd+xml" href="//en.wikisource.org/w/api.php?action=rsd"/>
23 <link rel="copyright" href="//creativecommons.org/licenses/by-sa/3.0/" />
24 <link rel="alternate" type="application/atom+xml" title="Wikisource Atom feed" href="/w/index.php?tit
25 <link rel="canonical" href="https://en.wikisource.org/wiki/George_Washington%27s_Fourth_State_of_the
26 <link rel="dns-prefetch" href="//login.wikimedia.org/" />
27 <link rel="dns-prefetch" href="//meta.wikimedia.org/" />
28 </head>
29 <body class="mediawiki ltr sitedir-ltr mw-hide-empty-elt ns-0 ns-subject page-George_Washington_s_F
30 <div id="mw-head-base" class="noprint"></div>
31 <div id="content" class="mw-body" role="main">
32 <a id="top"></a>
33
34 <div id="siteNotice" class="mw-body-content"><!-- CentralNotice --></div>
35 <div class="mw-indicators mw-body-content">
36 </div>
37 <h1 id="firstHeading" class="firstHeading" lang="en">George Washington's Fourth State of the
38 <div id="bodyContent" class="mw-body-content">
39 <div id="siteSub">From Wikisource</div>
40 <div id="contentSub"></div>
41 <div id="jump-to-nav" class="mw-jump">
```



Go from HTML-formatted text...



...to cleaner text ready for analysis.

```
1 [' Fellow-Citizens of the Senate and of the House of Representatives: ',
2 'It is some abatement of the satisfaction with which I meet you on the present occasion that, in felicitating you on a
3 continuance of the national prosperity generally, I am not able to add to it information that the Indian hostilities which
4 have for some time past distressed our Northwestern frontier have terminated.', 'You will, I am persuaded, learn with no less
5 concern than I communicate it that reiterated endeavors toward effecting a pacification have hitherto issued only in new and
6 outrageous proofs of persevering hostility on the part of the tribes with whom we are in contest. An earnest desire to procure
7 tranquillity to the frontier, to stop the further effusion of blood, to arrest the progress of expense, to forward the prevalent
8 wish of the nation for peace has led to strenuous efforts through various channels to accomplish these desirable purposes; in
9 making which efforts I consulted less my own anticipations of the event, or the scruples which some considerations were
10 calculated to inspire, than the wish to find the object attainable, or if not attainable, to ascertain unequivocally that
11 such is the case.', 'A detail of the measures which have been pursued and of their consequences, which will be laid before
12 you, while it will confirm to you the want of success thus far, will, I trust, evince that means as proper and as efficacious
13 as could have been devised have been employed. The issue of some of them, indeed, is still depending, but a favorable one,
14 though not to be despaired of, is not promised by anything that has yet happened.', 'In the course of the attempts which have
15 been made some valuable citizens have fallen victims to their zeal for the public service. A sanction commonly respected even
16 among savages has been found in this instance insufficient to protect from massacre the emissaries of peace. It will, I presume,
17 be duly considered whether the occasion does not call for an exercise of liberality toward the families of the deceased.',
18 'It must add to your concern to be informed that, besides the continuation of hostile appearances among the tribes north of
19 the Ohio, some threatening symptoms have of late been revived among some of those south of it.', 'A part of the Cherokees,
20 known by the name of Chickamaugas, inhabiting five villages on the Tennessee River, have long been in the practice of
21 committing depredations on the neighboring settlements.', 'It was hoped that the treaty of Holston, made with the Cherokee
22 Nation in July, 1791, would have prevented a repetition of such depredations; but the event has not answered this hope. The
23 Chickamaugas, aided by some banditti of another tribe in their vicinity, have recently perpetrated wanton and unprovoked
24 hostilities upon the citizens of the United States in that quarter. The information which has been received on this subject
25 will be laid before you. Hitherto defensive precautions only have been strictly enjoined and observed.', 'It is not understood
26 that any breach of treaty or aggression whatsoever on the part of the United States or their citizens is even alleged as a
27 pretext for the spirit of hostility in this quarter.', 'I have reason to believe that every practicable exertion has been made
28 (pursuant to the provision by law for that purpose) to be prepared for the alternative of a prosecution of the war in the event
29 of a failure of pacific overtures. A large proportion of the troops authorized to be raised have been recruited, though the
30 number is still incomplete, and pains have been taken to discipline and put them in condition for the particular kind of
```



Humanities data

- Data is material generated or collected while conducting research
- Examples of humanities data:
 - Citations
 - Code/Algorithms
 - Databases
 - Geospatial coordinates

Can you think of others?



Text as data

- Data quality
 - Clean vs. dirty OCR
 - HathiTrust OCR is dirty (uncorrected)
- Analyzed by corpus/corpora
 - Text corpus: a digital collection OR an individual's research text dataset
 - Text corpora: “bodies” of text
- Text decomposition/recomposition (Rockwell, 2003)
 - Cleaning data involves discarding data
 - Prepared text may be illegible to the human reader



Preparing data

A researcher may:

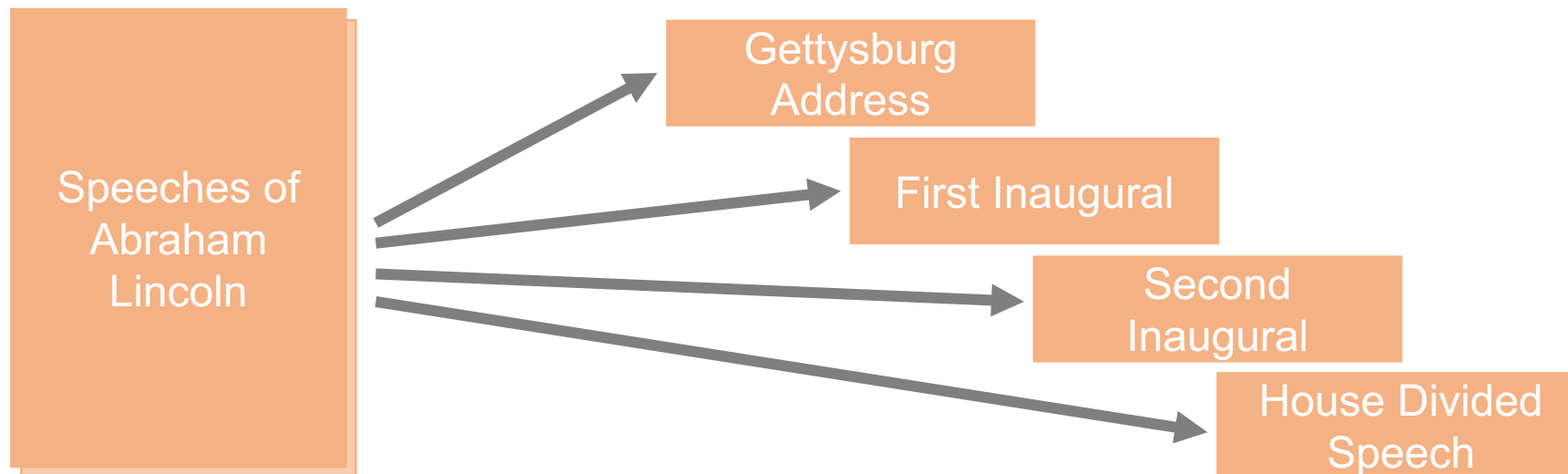
- Correct OCR errors
- Remove title, header information
- Remove html or xml tags
- Split or combine files
- Remove certain words, punctuation marks
- Lowercase text
- Tokenize the words



Key concepts

Chunking text

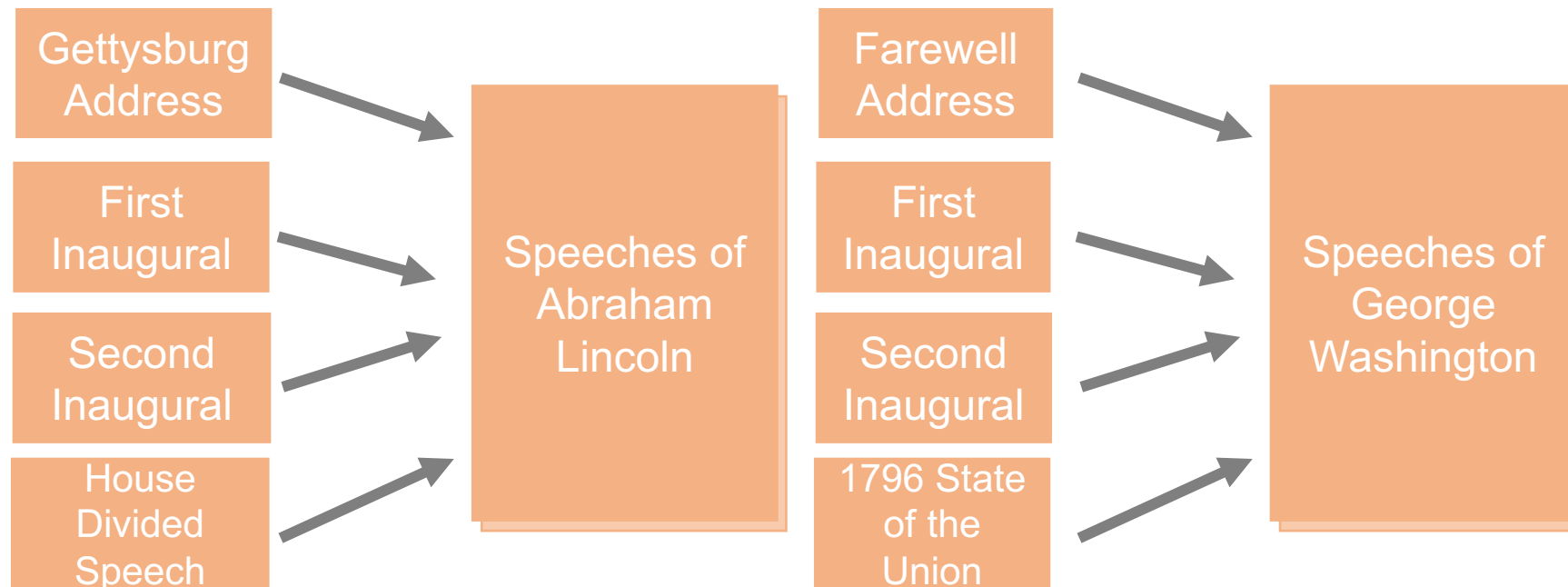
Splitting text into smaller pieces before analysis. May be divided by paragraph, chapter, or a chosen number of words (e.g. 1000 word chunks).



Key concepts

Grouping text

Combining text into larger pieces before analysis.



Key concepts

Tokenization

Breaking text into pieces called tokens. Often certain characters, such as punctuation marks, are discarded in the process

[four], [score], [and], [seven], [years], [ago], [our], [fathers], [brought], [forth], [on], [this], [continent], [a], [new], [nation], [conceived], [in], [liberty], [and], [dedicated], [to], [the], [proposition], [that], [all], [men], [are], [created], [equal]



Preparing data

- Preparation affects results
 - Amount of text and size of chunks
 - Which stop words removed; which characters are included
 - Whether to lowercase and normalize words
- Preparation for analysis takes time, effort
 - This is where scripting becomes useful!



Hands-on Activity

 *See Handout pp. 1-2*

- In groups of 2 or 3, assign each person several of the text preparation actions seen in the table to the right (Denny and Spirling, 2017).
- Read the descriptions. Then take turns explaining each to your group.

Term
Punctuation
Numbers
Lowercasing
Stemming
Stopword Removal
n-gram Inclusion
Infrequently Used Terms



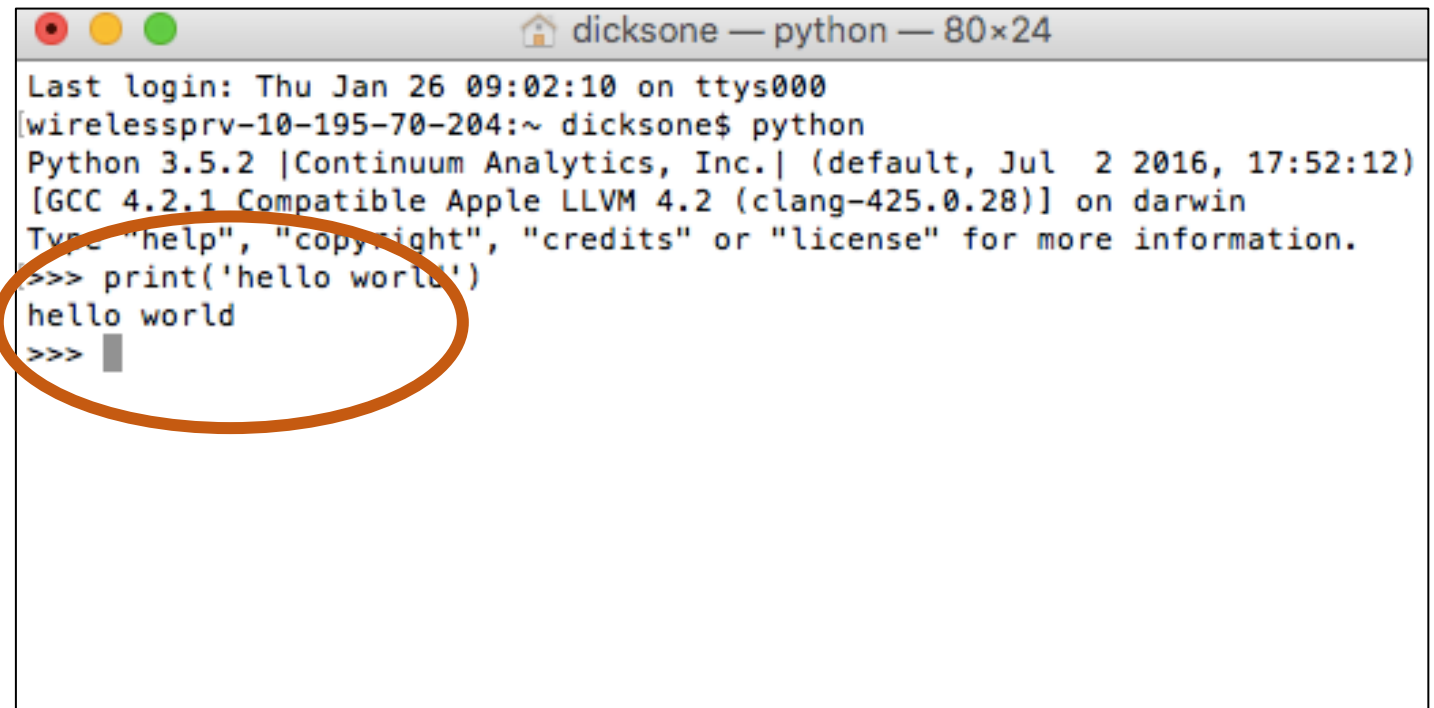
Scripting with Python

- Python is a scripting language
- Good for working with data
 - Interpreted language → follows step-by-step directions
- Relatively straightforward syntax
 - Avoids excess punctuation



Using Python: Interactive programming

- Using a **python interpreter**
- Run each step at the prompt
- If you enter “Python” on the command line, you will start the interpreter
- *We aren't using it today!*



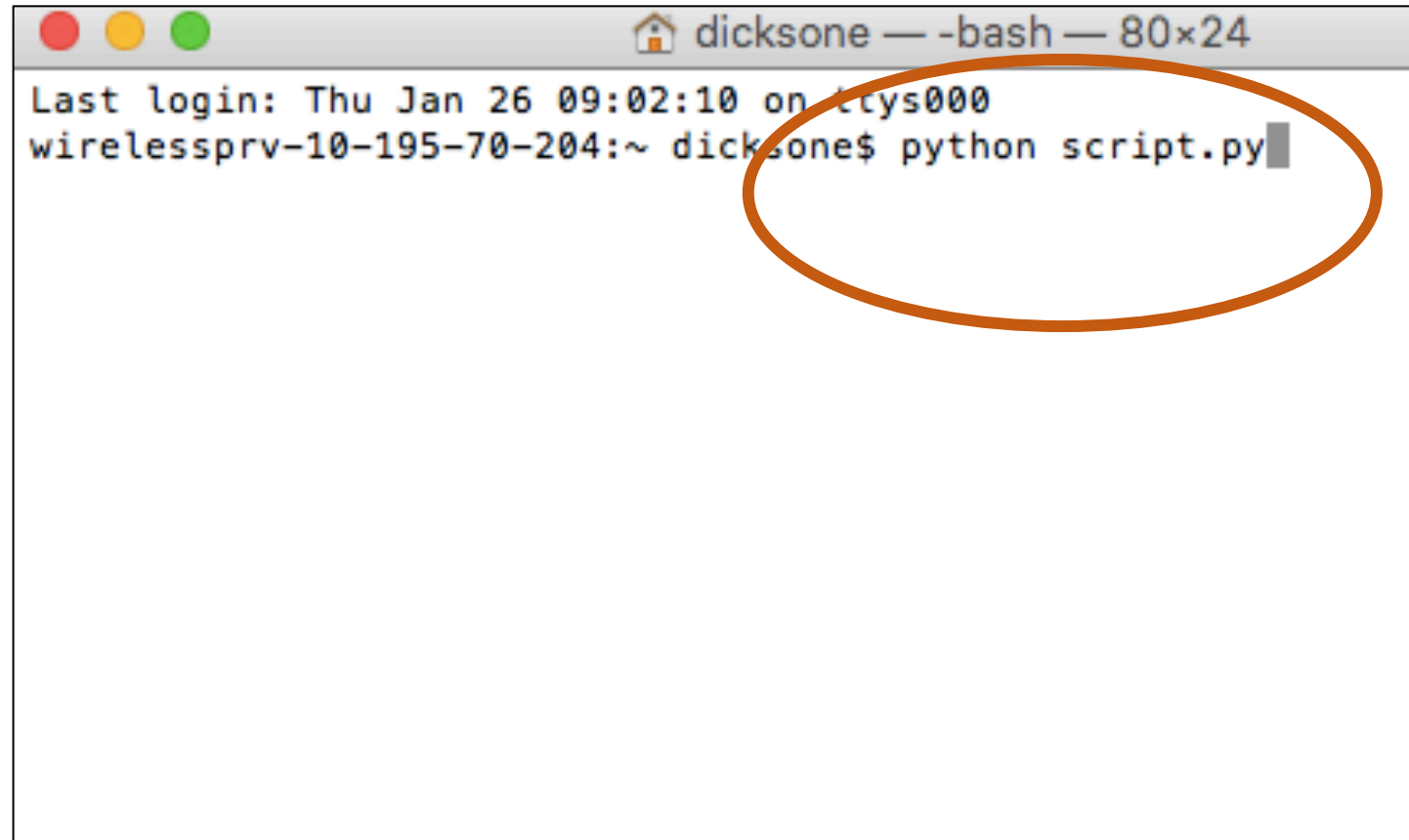
```
Home dickson — python — 80x24
Last login: Thu Jan 26 09:02:10 on ttys000
[wirelessprv-10-195-70-204:~ dickson]$ python
Python 3.5.2 |Continuum Analytics, Inc.| (default, Jul 2 2016, 17:52:12)
[GCC 4.2.1 Compatible Apple LLVM 4.2 (clang-425.0.28)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>> print('hello world')
hello world
>>> █
```



Using Python: Write & run scripts

- Scripts are directions for your computer to follow
- Save the script as a file ending in .py
- On the command line, run the script

→ This is how we'll do it!

A terminal window titled 'dicksone — -bash — 80x24' with a home icon. The window shows the following text: 'Last login: Thu Jan 26 09:02:10 on ttys000', 'wirelessprv-10-195-70-204:~ dicksone\$', and 'python script.py' with a cursor. The command 'python script.py' is circled in orange.

```
dicksone — -bash — 80x24
Last login: Thu Jan 26 09:02:10 on ttys000
wirelessprv-10-195-70-204:~ dicksone$ python script.py
```



Anatomy of running a Python script

```
python <script filename> <arguments>
```

1. Tells the computer to be ready for Python
2. Directs which file to run
3. The file may ask for additional information to be given at runtime, called arguments

Example: `python myscript.py newfile.txt`



Sample Reference Question

I'm a student in history who would like to incorporate digital methods into my research. I study American politics, and in particular I'd like to examine how concepts such as liberty change over time.

Approach:

Use a Python script to prepare the text data we scraped from Wikisource.



Hands-on activity

 See Handout p. 3

In this activity, you will run a Python script that will remove the HTML tags from George Washington's *Fourth State of the Union Address* you scraped earlier from Wikisource.



Strip HTML tags








































































- What You Need:
 - PythonAnywhere bash console
 - washington_4.txt
 - remove_tag.py script




Go to Files to view script

Files

Enter new file name, eg hello.py New file

 .bash_history	  	2017-03-16 06:40	802 bytes
 .bashrc	  	2016-03-28 15:42	559 bytes
 .gitconfig	  	2016-03-28 15:42	266 bytes
 .my.cnf	  	2017-04-13 14:12	40 bytes
 .mysql_history	  	2017-04-13 14:17	0 bytes
 .profile	  	2016-03-28 15:42	79 bytes
 .python_history	  	2017-02-09 18:20	56 bytes
 .pythonstartup.py	  	2016-07-01 20:53	77 bytes
 .sqlite_history	  	2017-04-14 19:53	2.1 KB
 .vimrc	  	2016-07-01 20:53	4.6 KB
 README.txt	  	2016-10-05 16:32	235 bytes
 get_verbs.py	  	2017-03-21 17:31	903 bytes
 mdp.49015002221837.json.bz2	 	2017-03-08 14:27	988.0 KB
 remove_stopwords.py	  	2017-05-05 14:37	353 bytes
 remove_tag.py	  	2017-07-07 13:17	374 bytes
 see_nouns_new.py	  	2017-07-05 15:58	590 bytes
 stopwords.txt	  	2017-05-05 14:37	621 bytes
 tagless_file.txt	  	2017-07-07 15:10	13.2 KB
 washington_4.txt	  	2017-06-22 19:00	41.1 KB

 Upload a file



View script

```
1 import sys
2 import re
3
4 file_name = sys.argv[1]
5
6 with open(file_name, 'r') as myfile:
7     contents = myfile.read().replace('\n', '')
8     pattern = r"<p>(.*?)</p>"
9     new_text = str(re.findall(pattern, contents))
10    new_text = new_text.replace("<i>", " ")
11    new_text = new_text.replace("</i>", " ")
12    new_text = new_text.replace("<a>", " ")
13    new_text = new_text.replace("</a>", " ")
14    new_text = new_text.replace("<b>", " ")
15    new_text = new_text.replace("</b>", " ")
16    cleaned = open('tagless_file.txt', 'w')
17    cleaned.write(new_text)
18
19
```



Open shell and run script












































































```
15:09 ~ $ python remove_tag.py washington_4.txt
```

Enter command:

```
python remove_tag.py washington_4.txt
```



Find new file

 .bash_history	   2017-03-16 06:40 802 bytes
 .bashrc	   2016-03-28 15:42 559 bytes
 .gitconfig	   2016-03-28 15:42 266 bytes
 .my.cnf	   2017-04-13 14:12 40 bytes
 .mysql_history	   2017-04-13 14:17 0 bytes
 .profile	   2016-03-28 15:42 79 bytes
 .python_history	   2017-02-09 18:20 56 bytes
 .pythonstartup.py	   2016-07-01 20:53 77 bytes
 .sqlite_history	   2017-04-14 19:53 2.1 KB
 .vimrc	   2016-07-01 20:53 4.6 KB
 README.txt	   2016-10-05 16:32 235 bytes
 get_verbs.py	   2017-03-21 17:31 903 bytes
 mdp.49015002221837.json.bz2	  2017-03-08 14:27 988.0 KB
 remove_stopwords.py	   2017-05-05 14:37 353 bytes
 remove_tag.py	   2017-07-07 13:17 374 bytes
 see_nouns_new.py	   2017-07-05 15:58 590 bytes
 stopwords.txt	   2017-05-05 14:37 621 bytes
 tagless_file.txt	   2017-07-07 15:10 13.2 KB
 washington_4.txt	   2017-06-22 19:00 41.1 KB



Review results

1 [' Fellow-Citizens of the Senate and of the House of Representatives: ',
2 'It is some abatement of the satisfaction with which I meet you on the present occasion that, in felicitating you on a
3 continuance of the national prosperity generally, I am not able to add to it information that the Indian hostilities which
4 have for some time past distressed our Northwestern frontier have terminated.', 'You will, I am persuaded, learn with no less
5 concern than I communicate it that reiterated endeavors toward effecting a pacification have hitherto issued only in new and
6 outrageous proofs of persevering hostility on the part of the tribes with whom we are in contest. An earnest desire to procure
7 tranquillity to the frontier, to stop the further effusion of blood, to arrest the progress of expense, to forward the prevalent
8 wish of the nation for peace has led to strenuous efforts through various channels to accomplish these desirable purposes; in
9 making which efforts I consulted less my own anticipations of the event, or the scruples which some considerations were
10 calculated to inspire, than the wish to find the object attainable, or if not attainable, to ascertain unequivocally that
11 such is the case.', 'A detail of the measures which have been pursued and of their consequences, which will be laid before
12 you, while it will confirm to you the want of success thus far, will, I trust, evince that means as proper and as efficacious
13 as could have been devised have been employed. The issue of some of them, indeed, is still depending, but a favorable one,
14 though not to be despaired of, is not promised by anything that has yet happened.', 'In the course of the attempts which have
15 been made some valuable citizens have fallen victims to their zeal for the public service. A sanction commonly respected even
16 among savages has been found in this instance insufficient to protect from massacre the emissaries of peace. It will, I presume,
17 be duly considered whether the occasion does not call for an exercise of liberality toward the families of the deceased.',
18 'It must add to your concern to be informed that, besides the continuation of hostile appearances among the tribes north of
19 the Ohio, some threatening symptoms have of late been revived among some of those south of it.', 'A part of the Cherokees,
20 known by the name of Chickamaugas, inhabiting five villages on the Tennessee River, have long been in the practice of
21 committing depredations on the neighboring settlements.', 'It was hoped that the treaty of Holston, made with the Cherokee
22 Nation in July, 1791, would have prevented a repetition of such depredations; but the event has not answered this hope. The
23 Chickamaugas, aided by some banditti of another tribe in their vicinity, have recently perpetrated wanton and unprovoked
24 hostilities upon the citizens of the United States in that quarter. The information which has been received on this subject
25 will be laid before you. Hitherto defensive precautions only have been strictly enjoined and observed.', 'It is not understood
26 that any breach of treaty or aggression whatsoever on the part of the United States or their citizens is even alleged as a
27 pretext for the spirit of hostility in this quarter.', 'I have reason to believe that every practicable exertion has been made
28 (pursuant to the provision by law for that purpose) to be prepared for the alternative of a prosecution of the war in the event
29 of a failure of pacific overtures. A large proportion of the troops authorized to be raised have been recruited, though the
30 number is still incomplete, and pains have been taken to discipline and put them in condition for the particular kind of



What happened

- Looks for `< >` to denote html tags
 - Using regular expressions and Python's `replace()`
- Writes the characters that are NOT html tags and their contents to a file
- Note: this script is simple, but not the most robust way to remove tags



On your own

 *See Handout p. 3*

Can you run a script to remove stop words?

- The script is called `remove_stopwords.py`
- Hint: The `remove_stopwords.py` script requires three arguments:
 - The input filename of your tagless file
 - The filename of the list of stop words
 - The output filename you make up
- Edit the `stopwords.txt` file to customize your list



Case Study: *Inside the Creativity Boom*

After downloading the Extracted Features data for the relevant volumes, used scripting to:

- Narrow corpus to individual pages that contained creativ*
 - Discarded all other pages
- Discard certain tokens such as pronouns and conjunctions
 - To keep only to most "meaningful" terms



Read and Reflect...

 *See Handout p. 4*

- Passage from “[Against Cleaning](#)” by Katie Rawson and Trevor Muñoz
- They suggest a strategy for dealing with humanities data:
 - Shared authority control across data sets
 - Indexes for nuance
 - Tidy, not clean data



Read and Reflect...

☞ *See Handout p. 4*

“When humanities scholars recoil at data-driven research, they are often responding to the reductiveness inherent in this form of scholarship. This reductiveness can feel intellectually impoverishing to scholars who have spent their careers working through particular kinds of historical and cultural complexity... From within this worldview, data cleaning is then maligned because it is understood as a step that inscribes a normative order by wiping away what is different. The term “cleaning” implies that a data set is ‘messy.’ “Messy” suggests an underlying order. It supposes things already have a rightful place, but they’re not in it—like socks on the bedroom floor rather than in the wardrobe or the laundry hamper.”

- “Against Cleaning” (Rawson and Muñoz, 2016)



Discussion

- *What does this excerpt suggest about the nuances of data cleaning?*
- *What does “clean” imply?*
- *How might you talk to researchers on your campus who would be uncomfortable with the idea of clean v. messy data?*



Questions?



References

- Denny, M. J. and Spirling, A. (2017). Text Preprocessing for Unsupervised Learning: Why It Matters, When It Misleads, and What to Do about It. <https://ssrn.com/abstract=2849145> .
- National Endowment for the Humanities. (2017) *Data Management Plans for NEH Office of Digital Humanities Proposals and Awards*. Retrieved October 1, 2017, from https://www.neh.gov/files/grants/data_management_plans_2018.pdf .
- Rawson, K., & Muñoz, T. (2016). Against Cleaning. Retrieved August 16, 2017, from <http://curatingmenus.org/articles/against-cleaning/> .
- Rockwell, G. (2003). What is Text Analysis, Really? *Literary and Linguistic Computing*, 18(2), 209–219. <https://doi.org/10.1093/lc/18.2.209> .

