

# “Smarter Data Collection for Microsoft Civic Graph

By Hannah Cutler

All of my work can be found in [this](#) GitHub repository.

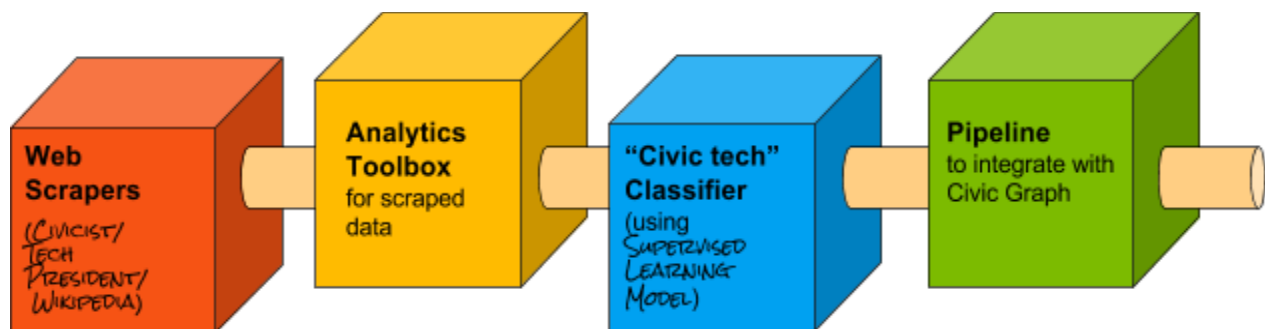
In this document:

- Overview
- Technical Documentation for: **Scrapers, Analysis, Classifier, Pipeline**
- Next Steps
- Thoughts and Ideas about: tools to consider using in future, incorporating additional external data sources into Civic Graph

## Overview

In an effort to make Civic Graph a little bit *smarter*, I developed **four** building blocks to help improve the quality (i.e. accuracy, completeness) of the data stored as well as automate aspects of the data collection process. They are:

1. **Web Scrapers**
2. **Analytics Toolbox**
3. **Classifier**
4. **Pipeline** to integrate scrapers, classifier, analysis with existing Civic Graph





## Technical Documentation

### Web Scrapers:

*I developed Python scripts to fetch all content from the Civicist and Tech President archives (published 2004-2015) as well as a script that pulls content from Wikipedia.*

Source	Relevant scripts	Libraries used	Notes/Status
Civicist Archives	<b>civ-link-parse.py:</b> <ul style="list-style-type: none"><li>- Parses Civicist Archives</li><li>- grabs all article URLs in Archive</li><li>- Writes article URLs to text file</li></ul> <hr/> <b>civ-article-parse.py:</b> <ul style="list-style-type: none"><li>- Read article URLs from text file</li><li>- Iterate through list of URLs and extract<ul style="list-style-type: none"><li>- Article title</li><li>- Article author</li><li>- Article date published</li><li>- Article content</li></ul></li><li>- Append all extracted content to string variable and write to text file</li></ul>	<ul style="list-style-type: none"><li>- <a href="#">BeautifulSoup</a></li></ul>	Active Accurate Not automated
TechPresident Archives	<b>tp-link-parse.py:</b> <ul style="list-style-type: none"><li>- Parses TechPresident Archives</li><li>- grabs all article URLs in Archive</li><li>- Writes article URLs to text file</li></ul> <hr/> <b>tp-article-parse.py:</b> <ul style="list-style-type: none"><li>- Read article URLs from text file</li><li>- Iterate through list of URLs and extract<ul style="list-style-type: none"><li>- Article title</li><li>- Article author</li><li>- Article date published</li><li>- Article content</li></ul></li><li>- Append all extracted content to string variable and write to text file</li></ul>	<ul style="list-style-type: none"><li>- BeautifulSoup</li></ul>	<i>TechPresident's Archives spanned 2004-2015. I discovered that over time the structure of the archive source code changed. There were 1411 pages total. I ran the script three times, extracting the article links in three batches.</i>
Wikipedia	<b>wikipedia_parser.py:</b> <ul style="list-style-type: none"><li>- Prompts user to enter a topic (and press ENTER)</li><li>- Fetches Wikipedia page contents for that topic if a page</li></ul>	<ul style="list-style-type: none"><li>- BeautifulSoup</li><li>- Python's wikipedia package</li></ul>	Active Accurate

	exists. Extracts: <ul style="list-style-type: none"> <li>- Page summary</li> <li>- Page Categories</li> <li>- Relevant links (pages associated with the topic)</li> <li>- Loops and continues to prompt for new query until user exits program</li> </ul>		
--	---	--	--

## Analysis (“Analytics Toolbox”)

*I experimented with several open source sentiment analysis/natural language processing tools. Based on the results of my experimentation, I decided to use spaCy and nltk.*

Tool/Library	Analysis performed	Relevant files	Notes/Status
AlchemyAPI	<ul style="list-style-type: none"> <li>- Entity Extraction (NER)</li> <li>- Keyword Extraction</li> <li>- Taxonomy extraction</li> <li>- Category Extraction</li> </ul>	<p><b>alchemyapi.py</b></p> <ul style="list-style-type: none"> <li>- Parses Civicist Archives</li> <li>- grabs all article URLs in Archive</li> <li>- Writes article URLs to text file</li> <li>- Run “python alchemyapi.py” and enter API key. Mine was: ff8f993db5ee0b907a3e41f19bbd57b8b4cbc24a</li> </ul> <hr/> <p><b>cg_categorize.py</b></p> <ul style="list-style-type: none"> <li>- Reads in file containing names of entities in Civic Graph and their associated descriptions</li> <li>- Attempts to categorize entities in Civic Graph based on descriptions using AlchemyAPI's "category" function</li> </ul> <hr/> <p><b>civictech_NER.py</b></p> <ul style="list-style-type: none"> <li>- Extracts data from Civicist/ TechPresident article content including:             <ul style="list-style-type: none"> <li>- Keywords</li> <li>- Taxonomy</li> </ul> </li> </ul>	<p><i>First tool that I used. After doing further analysis with other tools (i.e. nltk, spaCy), I learned that this tool is not very accurate, and it only extracted a fraction of the named entities. It was a useful exercise in teaching myself what kind of information one can extract from raw text using NLP libraries.</i></p> <p>Scripts retired.</p>

		<ul style="list-style-type: none"> <li>- Entities *</li> <li>- Categories</li> <li>- Function-</li> </ul>	
nlTK	<ul style="list-style-type: none"> <li>- Text tokenization</li> <li>- Entity Extraction (NER)</li> <li>- Concordance</li> <li>- Shared contexts for keywords</li> <li>- Frequency distribution</li> <li>- Collocation of common terms</li> </ul>	<p><b>nlTK_basic_analysis.py</b></p> <ul style="list-style-type: none"> <li>- Reads in raw text file and tokenizes text by sentence</li> <li>- Does semantic analysis on each token, extracting: <ul style="list-style-type: none"> <li>- All Named Entities</li> <li>- Unique entities</li> <li>- Concordance for a given keyword</li> <li>- Shared contexts for two or more words</li> <li>- Frequency distribution of named entities</li> <li>- collocations of commonly occurring terms</li> </ul> </li> </ul> <hr/> <p><b>reverse-madlibs.py</b></p> <ul style="list-style-type: none"> <li>- Reads in raw text file and tokenizes text by sentence</li> <li>- For each token: <ul style="list-style-type: none"> <li>- Extracts Named Entities</li> <li>- Computes all possible pairings between entities in that token and stores as a list of tuples</li> </ul> </li> <li>- Iterates through tuple list and prompts user to enter relationship between the two entities in a given tuple</li> <li>- Stores user response to external file</li> </ul>	<p><i>Accurate and easy to use but VERY slow on large datasets.</i></p> <p>Active</p>
<a href="#">spaCy</a>	<ul style="list-style-type: none"> <li>- Text tokenization</li> <li>- Part of speech (POS) labeling</li> <li>- Entity Extraction (NER)</li> <li>- Frequency distribution</li> </ul>	<p><b>spacy_experimentation.py</b></p> <ul style="list-style-type: none"> <li>- Reads in raw text file and tokenizes text by sentence using spaCy library</li> </ul> <hr/> <p><b>spacy_NER.py</b></p> <ul style="list-style-type: none"> <li>- Reads in raw text file and tokenizes text by sentence using spaCy library</li> </ul>	<p>Accurate and super fast.</p>

		<ul style="list-style-type: none"> <li>- For each token: <ul style="list-style-type: none"> <li>- Extracts Named Entities</li> <li>- Counts up entity frequencies and stores in dictionary {k: entity_name, v: freq}</li> <li>- Sorts dictionary in ascending order by value and prints resulting ordering</li> </ul> </li> </ul>	
General Analysis	Compute set intersection between two lists (e.g. Civic Graph entities and entities extracted from blog content)	<p><b>cg_queries.sql</b></p> <ul style="list-style-type: none"> <li>- SQL queries Civic Graph database to extract: <ul style="list-style-type: none"> <li>- Entity names, nicknames and id's from entity table</li> <li>- All connection types (collaboration, data, investment) from respective tables</li> </ul> </li> </ul> <hr/> <p><b>blog_graph_overlap.py</b></p> <ul style="list-style-type: none"> <li>- Reads in two csv files: <ul style="list-style-type: none"> <li>- entities pulled from Civic Graph database</li> <li>- Entities extracted from Civicist/ TechPresident blog content</li> </ul> </li> <li>- Computes set intersection between Civic Graph and blog entities</li> </ul> <hr/> <p><b>entities_analysis.R</b></p> <ul style="list-style-type: none"> <li>- Reads in two csv files: <ul style="list-style-type: none"> <li>- entities pulled from Civic Graph database</li> <li>- Entities extracted from Civicist/ TechPresident blog content</li> </ul> </li> <li>- Computes set intersection between Civic Graph and</li> </ul>	<b>blog_graph_overlap.py</b> and <b>entities_analysis.R</b> essentially do the same thing.

		<p>blog entities</p> <hr/> <p><b>entity_diff.py</b></p> <ul style="list-style-type: none"> <li>- Reads in two csv files: <ul style="list-style-type: none"> <li>- entities pulled from Civic Graph database</li> <li>- Entities extracted from Civicist/TechPresident blog content</li> </ul> </li> <li>- Computes set intersection between Civic Graph and blog entities</li> </ul>	
<a href="#">textacy</a>	N/A	Did not use this.	Did not use this.
<a href="#">Stanford nlp tools</a>	N/A	Did not use this.	Did not use this but apparently not as robust as nltk, textacy, or spaCy.

### Classifier:

1. **Input:** Raw text data
2. Tokenize training and test data
3. Label training data
4. Train classifiers on labeled training data
5. Test classifiers on test data
6. **Output:** For each token, returns an accuracy score for each label  $0.0 \leq \text{score} \leq 1.0$  corresponding to label relevance.

Functionality	Relevant scripts/functions	Notes/Status
Tokenize training data	<i>readfile_tokenizedata()</i> function within <b>algorithm.py</b>	Active
Labeling training data	<p><b>label_training_data.py</b></p> <ul style="list-style-type: none"> <li>- Read in each training dataset as text file</li> <li>- For each text file, generate labels for training data associated with a given label <ul style="list-style-type: none"> <li>- Funding</li> <li>- Data</li> <li>- Employment</li> <li>- Collaboration</li> <li>- Location</li> </ul> </li> </ul>	Active

Training classifiers	<code>svm()</code> function within <code>algorithm.py</code>	Active
Preprocess/ tokenize test data	<code>svm()</code> function within <code>algorithm.py</code>	Active <i>Tokenization on test data not fully implemented</i>
Run classifiers on test data	<code>accuracy()</code> function within <code>algorithm.py</code>	Active



## Pipeline:

View full-size diagram [here](#).

Editable version can be accessed [here](#).

## Next Steps (Task List)

### General tasks

- Cron job for all tasks within Mad Libs system.

### Scrapers

- Auto scraper for Civicist/TechPresident

### Classifier

- Choose an accuracy threshold (e.g.  $\geq 87\%$ )
- Make sure we are training the classifiers on the right data
- Structure the output of Classifier to fit schema of Civic Graph database
- Make sure we are training the classifiers on the right data
- Find the right source for our test data (is it Civicist/TechPresident? Something else?)
- Finish and tweak classifiers
- Talk to ML/NLP expert to get feedback (e.g. Jake Hofman)

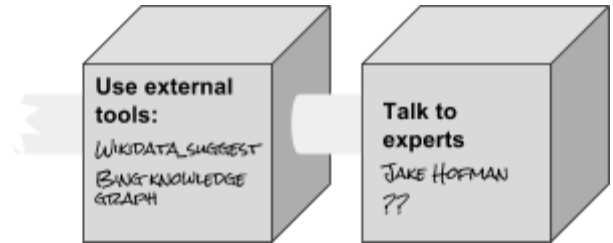
### Integration with Civic Graph Database

- Within **structure\_and\_post.py**, I wrote a pseudo-algorithm for:
  - structuring the output data from Classifier
  - Adding new table to database schema
  - POSTing structured data to Civic Graph database
  - *Note: Within this script, I describe all steps to integrate classified data with Civic Graph database.*

## Thoughts and Ideas

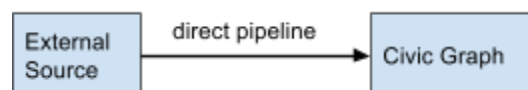
### Useful tools to consider in future:

- Some tools think could help us to gather additional data once we have a solid method for classifying civic tech relevance.
  - Run **wikipedia\_parser.py** on all “relevant” civic tech entities
  - Use Python’s [wikidata\\_suggest module](#) (To install: “pip install wikidata\_suggest”)
  - Use the Bing Knowledge Graph tool



### Suggestions for incorporating additional external data sources into Civic Graph:

- We have already identified potential sources that we can pull data from to incorporate into our Civic Graph database. These include:
  - DigitalSocialEU
  - Foundation Data
  - Civic Tech field guide
- As we expand this list, it will become increasingly necessary to think about how to utilize these resources to improve the quality of the data currently stored in civic graph (i.e. so that it is accurate and complete) and to think about how we can leverage these sources to identify new entities to add to Civic Graph. An obvious approach for incorporating external sources into the civic graph framework is to have direct pipelines from these sources to Civic Graph:



- You must assume that all of the data contained in the external source is relevant to civic tech and thus qualifies as being relevant to Civic Graph. Arguably, a better approach is to have an intermediate check where you aggregate and classify the external data before adding it to Civic Graph. For example, Create a “civic tech relevance score” and for each datapoint in the external source, if it generates a relevance score above some threshold, then pipe it to the Civic Graph database:

