

Natural Language processing in Python

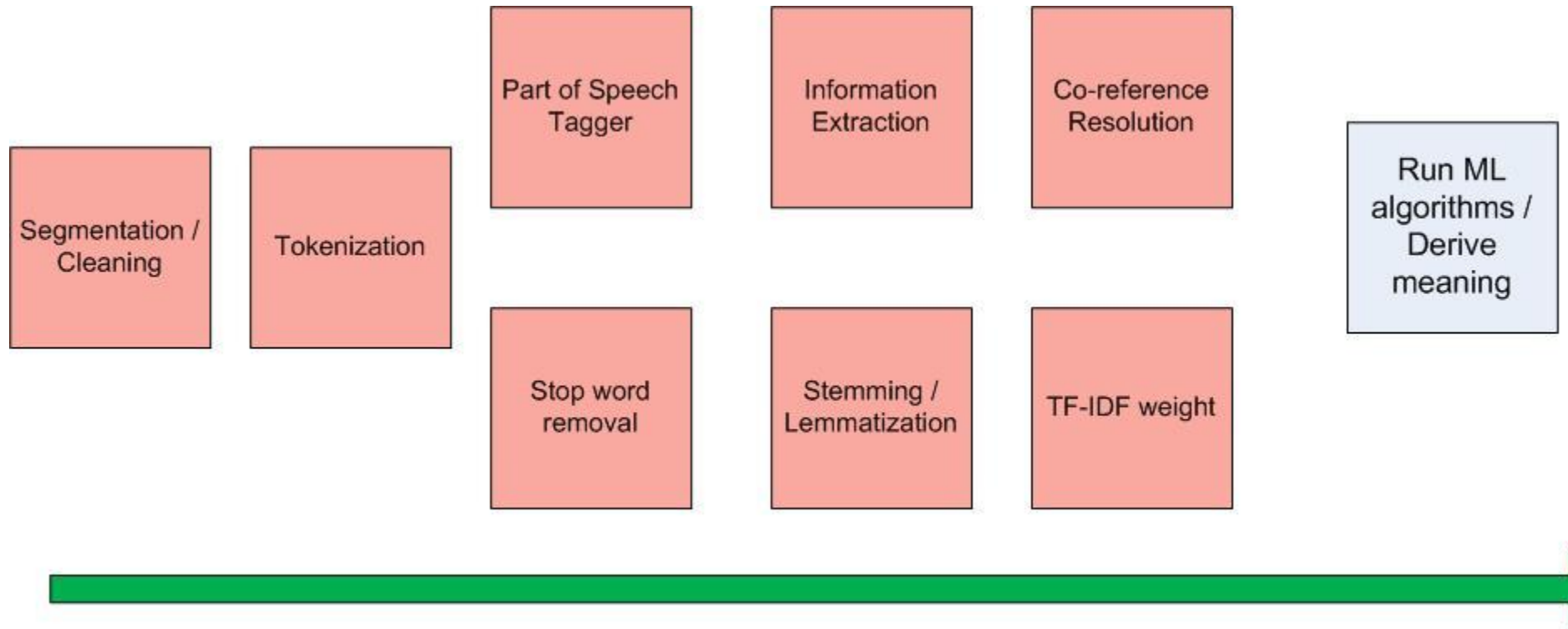
Chuck Chan

May 18, 2016

What is it used for

- Text classification
 - Document classification
 - information extraction
 - Sentiment analysis
- Automatic summarization
- Question answering

Pipeline



Python Libraries

- NLTK
 - Most comprehensive python library for NLP
 - Links to Stanford NER and POS Tagger
 - Apache 2.0 License
- Spacy
 - Cython based with OpenMP for multi-threaded generators
 - Used for speed and production
 - MIT License
- Gensim
 - Topic modeling, word2vec
 - Cython based
 - GNU Lesser General Public License

Cleaning text

- Remove html tags
 - Libraries: [BeautifulSoup](#)
- Deal with formatting issues
- Deal with encoding issues (Python 2.7)
 - Libraries: unicodedata, codecs
- No case insensitive

Pre-processing text

- Steps are to clean the data from the documents
 - Get rid of text that will bias classification algorithms
- Correctly weight words based on document other than by frequency
- Refine term counting so things are not double counted

Problems

- Word-sense disambiguation
 - Set[n]:a group or collection of things that belong together
 - Set[v] put, lay, or stand
 - Solved by relabeling with POS tagger and collocation dictionaries
- Metaphors and Similes
 - No reference point for comparisons
- Accuracy and speed

Tokenization

- Breaking a text body up into terms, symbols, and phrases.
- Tokenization can be used to separate out paragraphs, sentences, or terms.
- Usually separated by whitespace, sometimes include punctuation

Filtering Stop Words

- **Stop Word** a word that has the same likelihood of occurring in those documents not relevant to a query as in those documents relevant to the query.

Word	Count
the	3332
and	2972
a	1775
to	1725
of	1440
was	1161
it	1027
in	906
that	877
he	877

Top ten frequent words from Tom Sawyer

Types of words in stop words

- Determiners (the, a, an, another)
- Coordinating conjunctions (for, and, nor, but, or, yet, so)
- Prepositions (in, under, towards, before)
- Punctuation

List of stop words (Rank.nl)

• a	• did	• herself	• not	• the	• we've
• about	• didn't	• him	• of	• their	• were
• above	• do	• himself	• off	• theirs	• weren't
• after	• does	• his	• on	• them	• what
• again	• doesn't	• how	• once	• themselves	• what's
• against	• doing	• how's	• only	• then	• when
• all	• don't	• I	• or	• there	• when's
• am	• down	• I'd	• other	• there's	• where
• an	• during	• I'll	• ought	• these	• where's
• and	• each	• I'm	• our	• they	• which
• any	• few	• I've	• ours	• they'd	• while
• are	• for	• if	• ourselves	• they'll	• who
• aren't	• from	• in	• out	• they're	• who's
• as	• further	• into	• over	• they've	• whom
• at	• had	• is	• own	• this	• why
• be	• hadn't	• isn't	• same	• those	• why's
• because	• has	• it	• shan't	• through	• with
• been	• hasn't	• it's	• she	• to	• won't
• before	• have	• its	• she'd	• too	• would
• being	• haven't	• itself	• she'll	• under	• wouldn't
• below	• having	• let's	• she's	• until	• you
• between	• he	• me	• should	• up	• you'd
• both	• he'd	• more	• shouldn't	• very	• you'll
• but	• he'll	• most	• so	• was	• you're
• by	• he's	• mustn't	• some	• wasn't	• you've
• can't	• her	• my	• such	• we	• your
• cannot	• here	• myself	• than	• we'd	• yours
• could	• here's	• no	• that	• we'll	• yourself
• couldn't	• hers	• nor	• that's	• we're	• yourselves

Contractions & Compound words

- Problems
 - Examples: co-operation / cooperation may be counted as one or two tokens
 - Words such as San Francisco are counted twice for San and for Francisco
 - Possessives and contractions such as “Tom’s” and “you’re” may also be counted as two or more tokens

Compound terms

- Term types
 - Contractions: you're, we'll
 - Compound words: Los Altos, San Francisco, Facebook Inc.
 - Possessives: Tom's, Bill's
 - Collocations: **crystal clear, middle management**
- Identification
 - Using Named Entity Recognition(NER) to identify entities and tokenize them as one object
 - Chunking
 - Use bigrams to identify compound words manually
 - Have a dictionary or Regex
 - Use join function to combine with underscore
 - **Los Altos** → "Los_Altos" , you're → "you" "are"

Regular Expressions

- Search pattern mainly used for pattern matching
 - Library re
 - Cheat sheet: <http://regexlib.com/CheatSheet.aspx>
- Optimization tips
 - Consider speed vs. readability tradeoff.
 - Sparingly use and limit characters for greedy qualifiers [`*` or `+`]
 - Specify lengths `(hello){2}`
 - Use non capture group `(?:pattern)`
 - Extract common characters from alterations: use `th(is|at)` vs. `this|that`
 - Order alternations from most common to least `a(?:most|common|to|least)`
 - Use anchors for beginning`^` and end`$` of patterns

Special cases to consider

- Unicode characters
- Math equations, signs, numerical value
 - Fractions, decimals, and equations
- Twitter usernames, hashes, Emails addresses
 - #Hashtag, @Cruz
- Dates and Time
 - Previous dates, future dates, project time range
- Words with apostrophe or dashes
 - Contractions, hyphenated words
- Currency denominations
 - Large values, small values
- Chemical names, Genetic/protein sequences
- Account numbers
 - SSN, telephone, pins, userID

Reducing the number of distinct words

- Lemmatization
 - use of a vocabulary and morphological analysis of words
 - Slow and memory intensive but accurate

car, cars, car's, cars' => car

"better" has "good"

- Stemming
 - Use crude heuristic process that chops off the ends of words
 - Quick, but has false positives

; "universal", "university", and "universe" to "univers".

Part of speech tag

- Labels text to word classes
- Based on stochastic or rules based algorithms

Tag	Meaning	English Examples
ADJ	adjective	<i>new, good, high, special, big, local</i>
ADP	adposition	<i>on, of, at, with, by, into, under</i>
ADV	adverb	<i>really, already, still, early, now</i>
CONJ	conjunction	<i>and, or, but, if, while, although</i>
DET	determiner, article	<i>the, a, some, most, every, no, which</i>
NOUN	noun	<i>year, home, costs, time, Africa</i>
NUM	numeral	<i>twenty-four, fourth, 1991, 14:24</i>
PRT	particle	<i>at, on, out, over per, that, up, with</i>
PRON	pronoun	<i>he, their, her, its, my, I, us</i>
VERB	verb	<i>is, say, told, given, playing, would</i>
.	punctuation marks	<i>. , ; !</i>
X	other	<i>ersatz, esprit, dunno, gr8, univeristy</i>

Information Extraction

- Chunking & Chinking: Extraction of short phrases
 - Based on Part of speech patterns
 - Not very accurate
- NER: Task to extract and classify names of persons, organizations, locations, date time
 - Algorithm Conditional Random fields, Gazetteers, and rules

TF-IDF weights

- Reflect how important a word is to a document in a collection
- TF and IDF weights can have different weighing schemes depending on usage
- Weight increase with number of occurrence in a document and the rarity of the term in the collection

$$\text{tf-idf}_{t,d} = \text{tf}_{t,d} \times \text{idf}_t.$$

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \times \log_{10}(N / \text{df}_t)$$

Represent document as a weight vector of terms in vector space R

Term Occurrence	tf weight
0	0
1	1
2	1.3
10	2
1000	4

TF-IDF to vector space

Documents

		Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Terms	Antony	5.25	3.18	0	0	0	0.35
	Brutus	1.21	6.1	0	1	0	0
	Caesar	8.59	2.54	0	1.51	0.25	0
	Calpurnia	0	1.54	0	0	0	0
	Cleopatra	2.85	0	0	0	0	0
	mercy	1.51	0	1.9	0.12	5.25	0.88
	worser	1.37	0	0.11	4.15	0.25	1.95

Tf-idf weight for single word

Documents are represented as vector of weights

Term Frequency

- **Term frequency** measures how frequently a term occurs in a document
 - Generally represented as:
 - $Tf(t) = \frac{\text{(Number of times term } t \text{ appears in a document)}}{\text{(Total number of terms in the document)}}$

We use log scale because relevance does not increase proportionally with term frequency(count)

$$w_{t,d} = \begin{cases} 1 + \log_{10} tf_{t,d}, & \text{if } tf_{t,d} > 0 \\ 0, & \text{otherwise} \end{cases}$$

This offsets the case where a term occurs once ($\log_{10}(1) = 0$)

Inverse document Frequency

- **Inverse Document frequency** measure of how much information the word provides across all documents
 - measures how important a term is
 - $IDF(t) = \log \frac{\text{Total collection of documents}}{\text{Number of documents with term } t \text{ in it}}$
 - IDF (value of each word)
 - rare terms are more informative (should be high weight)
 - Frequent terms are less informative

$$\text{idf}_t = \log_{10} (N/\text{df}_t)$$

Log dampens score

Total
Number of documents

df = document frequency score: number of documents that contain term t, $\text{df} < N$

There is 1 IDF value for each term in a collection

Demonstration

**Latent Dirichlet allocation using
Gensim**