

Why Lemmatise?

- ◆ Lemmatisation is RollsRoyce stemming
- ◆ Only removes affixes if the resulting word is in its dictionary (so, is slower...)
- ◆ So, it should be more accurate; see Wordnet (Miller, 1995), Sketch Engine (Kilgarriff et al, 2004)
- ◆ Some systems can get quite complex, do partial parses and contain frequency info for lemmas

Why Stem?

REM

- ◆ So, we need to recognise variations of the same *stem* or *root**: fish~ for fishes, fishing, fished...
- ◆ We also need to recognise when two words are the same but from different syntactic categories (or parts of speech, POS); fish the *noun* and fish the *verb*
- ◆ Note, the root form of a word may be quite different be~ for is, are, am

* may be used differently

Lemmas: What's the Problem?

- ◆ It may get you better roots, than a stemmer (be for "am", "are", "is")
- ◆ Deals better with irregular plurals (eg woman and women)
- ◆ NB, can't itself recognise POS-based differences (fish and fish)

The screenshot shows the Wikipedia article on Lemmatisation. The page title is 'Lemmatisation'. The content discusses the process of grouping inflected forms of a word under a single lemma. It notes that many languages have several inflected forms, such as 'walk', 'walked', 'walks', 'walking'. The base form, 'walk', is called the lemma. The page also compares lemmatization with stemming, stating that lemmatization requires context and knowledge of the grammar of a language, while stemming is a simpler, rule-based process.

The screenshot shows a Python 3.4.1 shell window. The code imports the nltk library and uses the WordNetLemmatizer to attempt to lemmatize several words. The output shows that the word 'n' is not defined, resulting in a NameError. The full code and error message are as follows:

```
Python 3.4.1 (default, May 21 2014, 01:39:38)
[GCC 4.2.1 Compatible Apple LLVM 5.1 (clang-503.0.40)] on darwin
>>> import nltk
>>> wn = nltk.WordNetLemmatizer()
>>> wn.lemmatize('woman')
'woman'
>>> wn.lemmatize('women', 'v')
'woman'
>>> wn.lemmatize('fish', 'v')
'fish'
>>> wn.lemmatize('fishing', 'v')
'fishing'
>>> wn.lemmatize('is', 'v')
'be'
>>> wn.lemmatize('are', 'v')
'be'
>>> wn.lemmatize('am', 'v')
'be'
>>> n.lemmatize('am', 'n')
Traceback (most recent call last):
  File "<ipython-input-9>", line 1, in <module>
    n.lemmatize('am', 'n')
NameError: name 'n' is not defined
>>> wn.lemmatize('am', 'n')
'am'
>>>
```

Text Pre-Processing
Parts of Speech (POS)
& POS Tagging

Why POS Tag?

- We may also need to distinguish words that look the same but are from different syntactic categories (or parts of speech); *fish noun / verb*
- Lemmatising may need to know the part-of-speech already (noun or verb)
- Unfortunately, this may require parsing a whole sentence to disambiguate a POS and there are accuracy issues

POS Tagging Definition...

Create account Log in

Article Talk Read Edit View history Search

Part-of-speech tagging

From Wikipedia, the free encyclopedia

In corpus linguistics, **part-of-speech tagging** (POS tagging or POST), also called **grammatical tagging** or **word-category disambiguation**, is the process of marking up a word in a text (corpus) as corresponding to a particular part of speech, based on both its definition, as well as its context—i.e. relationship with adjacent and related words in a phrase, sentence, or paragraph. A simplified form of this is commonly taught to school-age children, in the identification of words as nouns, verbs, adjectives, adverbs, etc.

Once performed by hand, POS tagging is now done in the context of computational linguistics, using algorithms which associate discrete terms, as well as hidden parts of speech, in accordance with a set of descriptive tags. POS-tagging algorithms fall into two distinctive groups: rule-based and stochastic. E. Brill's tagger, one of the first and most widely used English POS-tagger, employs rule-based algorithms.

Penn Part of Speech Tags

Note: these are the 'modified' tags used for Penn tree banking; these are the tags used in the JET system. NP, NPS, PP, and PRPs from the original Penn part-of-speech tagging were changed to NNP, NNPS, PRP, and PRPS to avoid clashes with standard syntactic categories.

1. CC	Coordinating conjunction
2. CD	Cardinal number
3. DT	Determiner
4. EX	Existential there
5. FW	Foreign word
6. IN	Preposition or subordinating conjunction
7. JJ	Adjective
8. JJR	Adjective, comparative
9. JJS	Adjective, superlative
10. LS	List item marker
11. MD	Modal
12. NN	Noun, singular or mass
13. NNS	Noun, plural
14. NNP	Proper noun, singular
15. NNPS	Proper noun, plural
16. PDT	Pronominal determiner
17. POS	Possessive ending
18. PRP	Personal pronoun
19. PRP\$	Possessive pronoun
20. RB	Adverb
21. RBR	Adverb, comparative
22. RBS	Adverb, superlative
23. RP	Particle
24. SYM	Symbol
25. TO	to
26. UH	Interjection
27. VB	Verb, base form
28. VBD	Verb, past tense
29. VBG	Verb, present participle
30. VBN	Verb, past participle
31. VBP	Verb, non-3rd person singular present
32. VBZ	Verb, 3rd person singular present
33. WDT	Wh-determiner
34. WP	Wh-pronoun
35. WP\$	Possessive wh-pronoun
36. WRB	Wh-adverb

[Ln: 14 Col: 4]

Accuracy?

POS tagged tuples...

- There is a convention in Python to describe these tagged words in string tuples:

- 'fly/Vb', 'fly/NN', 'cheese/NN'

- So, you can write a sentence parse as a string that can be split and converted:

'The/DT man/NN sings/VB the/DT song/JJ'

[('The', 'DT'), ('man', 'NN'), ('sings', 'VB'), ('the', 'DT'), ('song', 'JJ')]

POS tagging...

- There are a plethora of parsers that can be used to recover the syntactic structure of sentences
- They can be used in conjunction with lemmatizers or part of them to get more accurate word identifications

```
Python 3.4.1 (default, May 21 2014, 01:39:38)
[GCC 4.2.1 Compatible Apple LLVM 5.1 (clang-503.0.40)] on darwin
Type "copyright", "credits" or "license()" for more information.

>>> import nltk
>>> text = nltk.word_tokenize("The fish jumped over the man who was fishing in the stream.")
[('The', 'DT'), ('fish', 'JJ'), ('jumped', 'VBD'), ('over', 'IN'), ('the', 'DT'), ('man', 'NN'), ('who', 'WP'), ('was', 'VBD'), ('fishing', 'VBG'), ('in', 'IN'), ('the', 'DT'), ('stream', 'NN'), ('.', '.')]
>>> text2 = nltk.word_tokenize("The fish sang the tune")
[('The', 'DT'), ('fish', 'JJ'), ('sang', 'NN'), ('the', 'DT'), ('tune', 'NN')]
>>> text3 = nltk.word_tokenize("The man sings the song")
[('The', 'DT'), ('man', 'NN'), ('sings', 'VNS'), ('the', 'DT'), ('song', 'JJ')]
>>> tag_str = '/The/DT man/NN sings/VB the/DT song/JJ'
>>> [nltk.tag.str2tuple(t) for t in tag_str.split()]
[('The', 'DT'), ('man', 'NN'), ('sings', 'VB'), ('the', 'DT'), ('song', 'JJ')]
>>> nltk.tag.str2tuple('man/NN')
('man', 'NN')
>>> tok = nltk.tag.str2tuple('man/NN')
>>> tok[1]
'NN'
>>> tok[0]
'man'
>>>
```

Why POS Tag?

- ◆ We may also need to distinguish words that look the same but are from different syntactic categories (or parts of speech); *fish noun / verb*
- ◆ Lemmatizing may need to know the part-of-speech already (noun or verb)
- ◆ Unfortunately, this may require parsing a whole sentence to disambiguate a POS and there are accuracy issues

REM

So, now we can lemmatise...

- ◆ POS tagging gives us an output we can submit to a lemmatizer
- ◆ However, note, the WordNet Lemmatizer generally works with simple tags ('n', 'v', 'adj') so you need to convert the more complicated penn-tags to use it

WordNet Lemmatizer

```
>>> import nltk
>>> text = nltk.word_tokenize('The fish who jumped over the man is happy, the man')
>>> text_with_pos = nltk.pos_tag(text)
>>> print(text_with_pos)
[('The', 'DT'), ('fish', 'NN'), ('who', 'NN'), ('jumped', 'VBD'), ('over', 'IN'),
 ('the', 'DT'), ('man', 'NN'), ('is', 'VBZ'), ('happy', 'JJ'), ('the', 'DT'), ('man', 'NN'),
 ('in', 'IN'), ('the', 'DT'), ('stream', 'NN')]
[('The', 'DT'), ('fish', 'NN'), ('who', 'NN'), ('jumped', 'VBD'), ('over', 'IN'),
 ('the', 'DT'), ('man', 'NN'), ('is', 'VBZ'), ('happy', 'JJ'), ('the', 'DT'), ('man', 'NN'),
 ('in', 'IN'), ('the', 'DT'), ('stream', 'NN')]

import nltk
text = nltk.word_tokenize('The fish who jumped over the man is happy, the man')
text_with_pos = nltk.pos_tag(text)
print(text_with_pos)

def convert_tags(tag):
    if tag == 'NN' or tag == 'VBD' or tag == 'VBZ':
        return 'n'
    else:
        return 'v'

wnl = nltk.WordNetLemmatizer()

for item in text_with_pos:
    new_tag = convert_tags(item[1].lower())
    print(item[0], new_tag)
    out = wnl.lemmatize(item[0], new_tag)
    print(out)
```

Old

Simpler Lemmatizer

```
>>> import nltk
>>> text = nltk.word_tokenize('John Doe ran the U.S., he\'ll do anything for I.B.M.')
>>> text_with_pos = nltk.pos_tag(text)
>>> print(text_with_pos)
[('John', 'NNP'), ('Doe', 'NNP'), ('ran', 'VBD'), ('the', 'DT'), ('U.S.', 'NNP'),
 ('.', '.'), ('he', 'PRP'), ('\'ll', 'MD'), ('do', 'VB'), ('anything', 'NN'), ('for',
 'IN'), ('I.B.M.', 'NNP')]

[('John', 'NNP'), ('Doe', 'NNP'), ('ran', 'VBD'), ('the', 'DT'), ('U.S.', 'NNP'),
 ('.', '.'), ('he', 'PRP'), ('\'ll', 'MD'), ('do', 'VB'), ('anything', 'NN'), ('for',
 'IN'), ('I.B.M.', 'NNP')]

import nltk
text = nltk.word_tokenize('John Doe ran the U.S., he\'ll do anything for I.B.M.')
text_with_pos = nltk.pos_tag(text)
print(text_with_pos)
ne_chunks = nltk.ne_chunk(text_with_pos, binary=True)
print(ne_chunks)
```

New

Text Pre-Processing

Parsing to Syntax

Of course,

- ◆ In general, the whole point of doing POS tagging and lemmatisation is to get to the syntactic structure of the sentence
- ◆ When you have the syntactic structure you can really separate out which bits are important (and disambiguate)
- ◆ **nltk** allows you to define grammars and use them (e.g., CFG = context-free grammar)

Focus on Parsing... REM

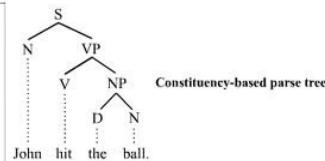
Parse - Merriam-Webster Online

www.merriam-webster.com/dictionary/parse
grammar : to divide (a sentence) into grammatical parts and identify the parts and their relations to each other. : to study (something) by looking at its parts ...

Parsing

Programming Language

Parsing or syntactic analysis is the process of analysing a string of symbols, either in natural language or in computer languages, according to the rules of a formal grammar. The term parsing comes from Latin pars, meaning part. [Wikipedia](#)



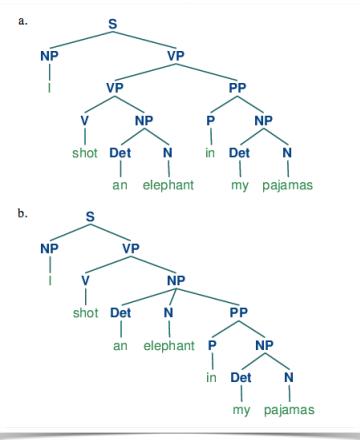
```

Python 3.4.1 (v3.4.1:1bc9e31e0105, Mar 18 2014, 00:54:21)
[GCC 4.8.1 (Apple Inc. build 5666) (dot 1)] on darwin
Type "copyright", "credits" or "license()" for more information.
>>> ===== RESTART =====
>>>
[s
  (NP I)
  (VP
    (VP (V shot) (NP (Det an) (N elephant)))
    (PP (P in) (NP (Det my) (N pajamas)))))
(s
  (NP I)
  (VP
    (V shot)
    (NP (Det an) (N elephant)) (PP (P in) (NP (Det my) (N pajamas))))))
>>>

import nltk
groucho_grammar = nltk.CFG.fromstring("""
S -> NP VP
NP -> P NP
NP -> Det N | Det N PP | 'I'
VP -> V NP
VP -> V PP
Det -> 'an' | 'my'
N -> 'elephant' | 'pajamas'
V -> 'shot'
P -> 'in'
PP -> 'in'
""")

sent = ['I', 'shot', 'an', 'elephant', 'in', 'my', 'pajamas']
parser = nltk.ChartParser(groucho_grammar)
for tree in parser.parse(sent):
    print(tree)
  
```

<http://www.nltk.org/book/ch08.html>



<http://www.nltk.org/book/ch08.html>

Stanford Parser

- The Stanford Parser is very commonly used to perform (better) parses of sentences
- Note, probabilistic parsers are often used because there are many alternative parses
- It is in Java but can be run from python wrappers (but, that is a story for another day)

Text Pre-Processing Spotting Entities

Our Focus... REM

- Text pre-processing is the poor-farmer cousin of full NLP; its not really about meaning
- Its about cleaning up text-data for future use
- Uses ideas from NLP (eg syntactic analysis, parsing) ... but is not often full NLP
- Ultimately, it seldom recovers meaning

Standard Tasks

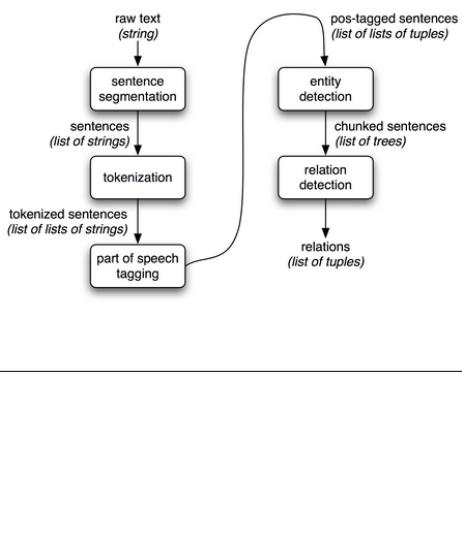
- ◆ *Tokenisation & Normalisation*: finding boundaries between word-like entities in character string
- ◆ *Fixing Misspellings*: where possible
- ◆ *Stemming, lemmatisation, POS-tagging*: finding slightly deeper identities between words (fished, fishing)
- ◆ *Removing Stop Words*: maximising the content-full words in the document/corpus
- ◆ *Entity Extraction*: identifying conceptual entities behind words

REM

Entities

- ◆ Entity extraction is the most semantic aspect of pre-processing (as such, often not done)
- ◆ Here, you identify the actual conceptual entity referred to by the word; encyclopaedic rather than dictionary knowledge
- ◆ I.B.M => "I.B.M" => **IBM** the organization

Typical Pipeline



Simple EG

```
Python 3.4.1 (v3.4.1:1bd3f2979cc4, May 18 2014, 00:54:21)
[GCC 4.2.1 (Apple Inc. build 5666) (dot 3)] on darwin
Type "copyright", "credits" or "license()" for more information.
>>>
[[('John', 'NNP'), ('Doe', 'NNP'), ('ran', 'VBD'), ('the', 'VBZ'), ('U.S.', 'NNP'), ('.', 'IN'), ('I.B.M', 'JJ'), ('.', 'IN')], [()])
(S
(NP John/NNP Doe/NNP)
ran/VBD
the/VBZ
(U.S./NNP)
.
I.B.M/JJ
.)
>>>
```

```
import nltk
text = nltk.word_tokenize('John Doe ran the U.S., he\'ll do anything for I.B.M.')
text with pos = nltk.pos_tag(text)
print(text with pos)
ne_chunks = nltk.ne_chunk(text with pos, binary=True)
print(ne_chunks)
```

Entity Extractors

- ◆ There are many different Entity Recognisers with different levels of accuracy for particular texts
- ◆ Stanford NER (Named Entity Recognizer)
- ◆ Open Calais is also heavily used
- ◆ But, need good reasons to go this far in your pre-processing, esp. considering accuracy

Text Pre-Processing Removing Stop Words

Why Remove Stops?

- ❖ We have been mainly transforming words in various ways to make them better for use...
- ❖ But, some words do not help, like stop words
- ❖ They are words that do not convey much content, they are not *contentful*
- ❖ They are also often very frequent and can effects norms and counting (cf Lect4)

Stops Words...(lucene)

"a", "an", "and", "are", "as", "at", "be", "but", "by", "for", "if", "in", "into", "is", "it", "no", "not", "of", "on", "or", "such", "that", "the", "their", "then", "there", "these", "they", "this", "to", "was", "will", "with"

Another set of stop-words

a,able,about,across,after,all,almost,also,am,among,an,an
d,any,are,as,at,be,because,been,but,by,can,cannot,coul
d,dear,did,do,does,either,else,ever,every,for,from,get,go
t,had,has,have,he,her,hers,him,his,how,however,i,if,in,i
nto,is,it,its,just,least,let,like,likely,may,me,might,most,
must,my,neither,no,nor,not,of,off,often,on,only,or,other,
our,own,rather,said,say,says,she,should,since,so,some,t
han,that,the,their,them,then,there,these,they,this,tis,to,t
oo,twas,us,wants,was,we,were,what,when,where,whic
h,while,who,whom,why,will,with,would,yet,you,your

Stop Word Lists...

- ❖ There is no definitive stop-word set, may vary for different purposes (see wikipedia for some egs)
- ❖ They can result in large reductions (40%-60%) in normal texts, so you are left with core content...

```
Python 3.4.1 (default, May 21 2014, 01:39:38)
[GCC 4.2.1 Compatible Apple LLVM 5.1 (clang-503.0.40)] on darwin
Type "copyright", "credits" or "license()" for more information.
>>> import nltk
>>> from nltk.corpus import stopwords
>>> stop = stopwords.words('english')
>>> stop
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', 'your', 'you
rs', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'her', 'hers',
'herself', 'it', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves',
'what', 'which', 'who', 'whom', 'this', 'that', 'these', 'those', 'am', 'is', 'are',
'was', 'were', 'been', 'being', 'have', 'has', 'had', 'having', 'd', 'do', 'does',
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as',
'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between',
'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to',
'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further',
'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any',
'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such', 'no', 'nor', 'not',
'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will',
'just', 'don', 'should', 'now']
```

```
>>> text = open('/Users/user/Desktop/text.txt')
>>> rawtext = text.read()
>>> [i for i in rawtext.split() if i not in stop]
['So', 'bunch', 'text', 'i've', 'put', 'together', 'check', 'tokenisation', 'crap',
'I', 'interested', 'I.B.M.', 'U.S.A.', 'USA', 'handled', 'ok?', 'This', 'what',
her, 'properly', 'deals', 'websites', 'like', 'www.ucd.ie', 'email', 'addresses',
'like', 'mark.keane@ucd.ie.', 'Also', 'weirdities', 'like', 'great', 'O'N
eill', 'like', 'M*A*S*H', 'maybe', '7', 'tokens?']
```

```
>>> rawtext
"So, this is just a bunch of text that i've put together to check on\ nthis tokenis
ation crap and I am interested in how I.B.M. or the U.S.A. or the USA\ s handled
, ok? This and whether it properly deals with websites like www.ucd.ie\ nand email
addresses like mark.keane@ucd.ie. Also other weirdities like\ the great O'
Neill and like M*A*S*H, which should be maybe 7 tokens?\n"
>>>
```

Pre-Processing

When to Use What & When...

When to use these things...

- ◆ Pre-processing is an important design choice: sometimes you might only do stemming and stop-word removal; some times lemmatization and stop-word removal
- ◆ Other times, you want to leave everything in and do full parsing of the text; or parsing and then stop-word removal later

Pre-Processing for Indexing...

- ◆ Indexing and retrieval tend to use quite crude stemming and stop-word removal (cf Lucene)
- ◆ Here you want a small no of string-features to capture your docs and queries; can increase recall without damaging precision
- ◆ Does not matter if they are crap (lying -> li) because they are consistent over the doc set and the scale of the corpus irons out inaccuracies

Pre-Processing for Meaning...

- ◆ Stemmers do not often get at morphological root; so, less good if you need more meaning
- ◆ Cases needing meaning invite lemmatization (POS and stop-word removal)
- ◆ Eg if you want definite verbs and nouns and other POS (cf Gerow & Keane, 2011)

Pre-Processsing for Tweets...

- ◆ In Twitter, eg, everything is different:
 - ◆ Twitter-specific pos-tagger
 - ◆ Stop-word removal can damage performance
 - ◆ Misspellings and abbreviations need specific treatment

Some Twitter Refs...

Gimpel, K., Schneider, N., O'Connor, B., Das, D., Mills, D., Eisenstein, J., ... & Smith, N. A. (2011, June). Part-of-speech tagging for twitter. In COLING: Volume 2 (pp. 42-47). ACL

Kouloumpis, E., Wilson, T., & Moore, J. (2011). Twitter sentiment analysis: The good the bad and the omg!. ICWSM, 11, 538-541.

Aiello, L. M., Petkos, G., Martin, C., Corney, D., Papadopoulos, S., Skraba, R., ... & Jaimes, A. (2013). Sensing trending topics in Twitter. IEEE Trans Multimedia.

Handling Different Sources...

Different file formats...

- ◆ Text files (we've seen)
- ◆ Web pages, html and xml
- ◆ PDFs and Docs

We have seen text.file input

A screenshot of a Mac desktop. In the top right corner, there's a system status bar showing icons for battery, signal, and time (Tue 15:43). Below it is a dock with several icons: 'Aktivit鋞en', 'MobileMac', 'Proj Programm', 'X. All Readings', 'X. Research', 'X. Teaching', 'X Admin', and 'X Tools'. The main screen shows a Python 3.4.1 terminal window titled 'Python 3.4.1 Shell'. It contains code demonstrating file reading and tokenization. To the right of the terminal is a text editor window titled 'Python 3.4.1: text.txt - /Users/user/Desktop/text.txt' showing the same text as the terminal. At the bottom of the screen is a file browser window showing various PDF files.

```
Python 3.4.1 (default, May 21 2014, 01:39:38)
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-503.0.40)] on darwin
Type "copyright", "credits" or "license()" for more information.
>>> ffile = open('/Users/user/Desktop/text.txt')
>>> rawtext = ffile.read()
>>> tokens = re.compile(r'\w+').findall(rawtext)
>>> tokens
['this', 'is', 'just', 'a', 'bunch', 'of', 'text', 'that', 'is', 'you',
 'put', 'together', 'to', 'check', 'on', 'this', 'tokenisation', 'crap', 'and', 'i',
 'am', 'interested', 'in', 'how', 'I.B.M.', 'or', 'the', 'U.S.A.', 'is', 'ha-
 nded', 'ok', 'when', 'dealing', 'with', 'websites', 'like', 'www.ucd.ie', 'and', 'email', 'addresses', 'like', 'mark.keane@ucd.
 ie', 'and', 'like', 'Also', 'other', 'weirdities', 'like', 'the', 'great',
 'o'Neill', 'and', 'like', 'www-e-g.com', 'which', 'should', 'be', 'maybe', '?',
 'tokens', '?']
```

We have seen text.file input

A screenshot of a Mac desktop. In the top right corner, there's a system status bar showing icons for battery, signal, and time (Mon 19:12). Below it is a dock with icons for Finder, Mail, Safari, and others. The main screen shows a Python 3.4.1 terminal window titled 'Python 3.4.1 Shell'. It contains code demonstrating file reading and manipulation. To the right of the terminal is a text editor window titled 'Python 3.4.1: simple.txt - /Users/user/Desktop/simple.txt' showing the string 'simple string to read in to python.'. At the bottom of the screen is a file browser window showing various PDF files.

```
Python 3.4.1 (default, May 21 2014, 01:39:38)
[GCC 4.2.1 Compatible Apple LLVM 5.0 (clang-503.0.40)] on darwin
Type "copyright", "credits" or "license()" for more information.
>>> open('./Users/user/Desktop/simple.txt').read()
<string 'simple string to read in to python.'>
>>> open('./Users/user/Desktop/simple.txt').read()
<string 'simple string to read in to python.'>
>>> string = open('./Users/user/Desktop/simple.txt').read()
<string 'simple string to read in to python.'>
>>> type(string)
<type 'str'>
>>> 
```

We have seen text.file input

- ◆ You basically open the file; `open()`
- ◆ Read it in: `read()`
- ◆ Then you have a string object you can manipulate in different ways

Web Pages are trickier...

- ◆ You basically open the file with a special method for opening ; `openurl.request.url()`
- ◆ Read it in: `read()`
- ◆ Then you need a new package (`BeautifulSoup`) to create an new object you can extract different bits of the page from

Installation steps...

- ◆ NB. Stuff in book is legacy; `clean_html` and `urlopen` do not work as specified
- ◆ Install correct version of BeautifulSoup for Python3.4 (aaaagghhh...called bs4 !)
<http://www.crummy.com/software/BeautifulSoup/>
- ◆ The import it and use its given commands

Reading Webpages...

```
Python 3.4.1 (default, May 21 2014, 01:39:38)
[GCC 4.2.1 Compatible Apple LLVM 5.1 (clang-503.0.40)] on darwin
Type "copyright", "credits" or "license()" for more information.
>>> import urllib
>>> url = 'http://www.csi.ucd.ie/users/mark-keane'
>>> urllib.request.urlopen(url)
Traceback (most recent call last):
  File "<ipython-input-2>", line 1, in <module>
    urllib.request.urlopen(url)
AttributeError: 'module' object has no attribute 'request'
>>> import urllib.request
>>> urllib.request.urlopen(url)
<http.client.HTTPResponse object at 0x10fcf46a0>
>>> rawhtml = urllib.request.urlopen(url).read()
>>> from bs4 import BeautifulSoup
>>> soup = BeautifulSoup(rawhtml)
>>> soup = bs4.BeautifulSoup(rawhtml)
Traceback (most recent call last):
  File "<ipython-input-3>", line 1, in <module>
    File "", line 1, in <module>
      File "<bs4>.BeautifulSoup(rawhtml)
NameError: name 'bs4' is not defined
>>> import bs4
>>> soup = bs4.BeautifulSoup(rawhtml)
>>> soup.title
<title>Mark Keane | UCD School of Computer Science and Informatics</title>
>>> soup.body.get_text(strip=True)
'UCD School of Computer Science and InformaticsScoil na Riomheolaiochta agus na F
aistíneacha UCDCalendarNewsPeopleSite mapSign In You are here:Home>CSI People>Ma
rk KeaneBiographyResearch Interests:AnalogyCognitive ScienceEvolutionSI
milarityText AnalyticsGeneralName and Title:Professor Mark Keane BA MA PhDPosition
n:Chair of Computer SciencePhone:Ext. 2470Email:Office:CSI / B2.01Address:School
of Computer Science & InformaticsBiographyProfessionalPublicationsResearchBiograp
hysince 1998, Prof. Mark Keane has been Chair of Computer Science at University C
```

Parsing Webpages...

```
>>> foollinks = soup.findAll('a')
>>> for link in foollinks: print(link)

<a id="navigation-top" name="top"></a>
<a href="/" rgl="home" title="Home"></a>
<a href="/" rgl="home" title="Home">Home</a>
  UCD School of Computer Science and Informatics      </a>
<a class="menu-1-1" href="#">Home > Calendar</a>
<a class="menu-1-2" href="#">News > News</a>
<a class="menu-1-3" href="#">Content > People</a>
<a class="menu-1-4" href="#">Sitemap > Display a site map with RSS feeds</a>
<a class="menu-1-5" href="#">Site map > Site map</a>
<a class="menu-1-6" href="#">Sign in > Sign in</a>
<a href="#">Home</a>
<a href="#">CSI People</a>
<a class="taxonomy_term_346" href="#">Category > Research-Interests > Analogy</a>
<a class="taxonomy_term_181" href="#">Category > Research-Interests > Cognitive Science</a>
<a class="taxonomy_term_949" href="#">Category > Research-Interests > Evolution</a>
<a class="taxonomy_term_378" href="#">Category > Research-Interests > Similarity</a>
<a class="taxonomy_term_940" href="#">Category > Research-Interests > Text Analytics</a>
<a href="#">Tabs-Tabset-1 > Biography</a>
<a href="#">Tabs-Tabset-2 > Professional</a>
<a href="#">Tabs-Tabset-4 > Research</a>
<a href="https://rms.ucd.ie/ufrs/W_RMS_PUB_COMMON_PUB_POPUP?p_object_id=134951287" target="_blank">> Details</a>
<a href="https://rms.ucd.ie/ufrs/W_RMS_PUB_COMMON_PUB_POPUP?p_object_id=8272699" target="_blank">> Details</a>
<a href="https://rms.ucd.ie/ufrs/W_RMS_PUB_COMMON_PUB_POPUP?p_object_id=82857200" target="_blank">> Details</a>
<a href="https://rms.ucd.ie/ufrs/W_RMS_PUB_COMMON_PUB_POPUP?p_object_id=88772700" target="_blank">> Details</a>
<a href="https://rms.ucd.ie/ufrs/W_RMS_PUB_COMMON_PUB_POPUP?p_object_id=134951303" target="_blank">> Details</a>
<a href="https://rms.ucd.ie/ufrs/W_RMS_PUB_COMMON_PUB_POPUP?p_object_id=134951297" target="_blank">> Details</a>
```

PDFs & .Docs...

- ◆ General advice is to convert them to text and work from there...
- ◆ But PDFminer is a python package for parsing PDFs (a bit complicated)

Selecting a Corpus

Finally, we have assumed...

- ◆ That you just know which texts to pre-process; but, sometimes you have to think about selecting the texts that make up a corpus
- ◆ Is this defined naturally; every debate in the Dail since 1922... (simple case)
- ◆ Every news article about stock markets... how do we define this? (medium case)
- ◆ Every tweet that is about senate elections ... how do we define this? (hard case)

Please Go Home Now...

- ◆ At least, if you have finished the practical...