

N-grams in Python

L445/L515
Autumn 2011

Calculating n-grams

We want to take a practical task, i.e., using n-grams for natural language processing, and see how we can start implementing it in Python.

Some of today will move rather fast

- You may not have been able to come up with the final program on your own ...
- But I want you to be able to understand how it works

2

Defining the problem

Here's what I want us to do today:

- Input a text file
- Output a file containing each bigram from the text, together with its frequency.

3

Skeleton of the Inputting

```
# Declare file to be worked with
textfile = "furniture.txt"

# Open & read file
for line in open(textfile):
    line = line.rstrip()

    print line
```

4

What Next?

Okay, we're able to read in a file line-by-line. Now, we need to do something ...

1. Tokenize each line into words
2. Create bigrams from the individual words
3. Store the bigrams

For the moment, let's pretend we're working with unigrams, and so we'll do the following:

1. Tokenize each line into words
2. Store the unigrams

5

Tokenization

A simple way to tokenize everything is to use `split`, which takes a string, splits it on whitespaces, and returns every non-whitespace item in a list.

```
for line in open(textfile):
    line = line.rstrip()

    # tokenize the text:
    tokens = line.split()
    ...
    # do stuff here
    ...
```

6

Side topic: modules

That's a pretty simplistic tokenizer. (Maybe someone's written a better one at some point ...)

In fact, I wrote a slightly better tokenizer and put it in a file called `useful.py`

- Q: How can you use my better tokenizer?
- A: With the python `import` statement, which allows you to import *modules*

7

import

At the top of your file, include the line:

- `from useful import tokenize`
- This says: from the module `useful` import the function `tokenize`

And then when we need the tokenizer, we can call it:

```
for line in open(textfile):
    line = line.rstrip()
    # tokenize the text:
    tokens = tokenize(line)
    # do stuff here
```

8

Unigrams

Thus far, we've: read in a file and tokenized each line. Thus, we have access to unigrams.

```
for line in open(textfile):
    line = line.rstrip()
    # tokenize the text:
    tokens = tokenize(line)

    # loop over unigrams:
    for word in tokens:
        print word
```

9

Data storage

Great! We have unigrams. ... And now we want to store them somewhere, along with their associated frequencies

- A list isn't quite right: once we store a word, how do we access it in the list when we encounter it again in a text?
- Plus, the data structure we want should store each word (type) and have the word "point to" its associated count

Python dictionaries are what we want ...

10

Dictionaries

Python has a **dictionary** data type, which is kind of like a set (unordered, no repeat keys), but is technically a hash which maps from one item (key) to another (value)

- It is an unordered set of *key:value* pairs

```
>>> d1 = {}                                # empty dictionary
>>> d1 = {'a':1, 'b':2, 'c':3}             # dictionary with 3 keys
>>> x = d1['a']                             # x = 1
>>> del d1['c']                             # d1 = {'a':1, 'b':2}
>>> d1['d'] = 4                             # d1 = {'a':1, 'd':4, 'b':2}
>>> k = d.keys()                           # k = ['a','d','b']
```

11

More on dictionaries

- To test whether an item is in a dictionary, use: `'a' in d1`
- To iterate over the items in a dictionary, you have lots of options:

```
for alpha, num in d1.iteritems():
    print alpha + '\t' + str(num)
for alpha in d1.keys():
    print alpha + '\t' + str(d1[alpha])
```

Most commonly, I use the following:

```
for alpha in d1:
    print alpha + '\t' + str(d1[alpha])
```

12

Using a dictionary to store unigrams

So, let's create a dictionary `Unigrams`, which we'll use to store each unigram

- Creation of dictionary: `Unigrams = {}`
- Adding each word to the dictionary: `Unigrams[word] = 1` (not quite right)

```
Unigrams = {}
for line in open(textfile):
    line = line.rstrip()
    # tokenize the text:
    tokens = tokenize(line)

    # loop over unigrams:
    for word in tokens:
        Unigrams[word] = 1
```

13

Adding each unigram

We really want to do the following:

- Make `Unigrams[word] = 1` if we've never seen this word before
- Make `Unigrams[word] = Unigrams[word] + 1` if we have seen this before
 - Use an if statement!

```
# loop over unigrams:
for word in tokens:
    if word in Unigrams:
        Unigrams[word] += 1
    else:
        Unigrams[word] = 1
```

14

Outputting unigrams to a file

- Iterate over the unigrams in the dictionary
- Write out each unigram, along with its count

```
# Write unigrams to output file:
output_file = open('unigrams.txt', 'w')
for unigram in Unigrams:
    count = Unigrams[unigram]
    output_file.write(str(count)+'\t'+unigram+'\n')
output_file.close()
```

15

From unigrams to bigrams

Now, let's extend the analysis from unigrams to bigrams.

- What do we need to do to make this happen? (Aside from nicely renaming our dictionary to `Bigrams` and so forth).

```
# loop over words in input:
for word in tokens:
    # something more needs to happen here:
    bigram = word
    if bigram in Bigrams:
        Bigrams[bigram] += 1
    else:
        Bigrams[bigram] = 1
```

16

Calculating bigrams

One thing novices sometimes do is to tokenize each line into tokens and then loop over tokens from 1 to $i-1$, taking pairs of words from the list

- Problem: doesn't account for bigrams between lines

Better solution: Keep track of the previous word

- Before we start looping over the input, include these lines:

```
# initialize variables:
Bigrams = {}
prev_word = "START"
```

17

Calculating bigrams (cont.)

```
# loop over words in input:
for word in tokens:
    # concatenate words to get bigram:
    bigram = prev_word + ' ' + word
    if bigram in Bigrams:
        Bigrams[bigram] += 1
    else:
        Bigrams[bigram] = 1
    # change value of prev_word
    prev_word = word
```

18

<div data-bbox="159 178 651 212" data-label="Section-Header"><h2>Regular expressions (for better tokenizing)</h2></div> <div data-bbox="110 279 699 306" data-label="Text"><p>[If we have time, we'll cover these slides ... but we probably won't]</p></div> <div data-bbox="27 325 734 352" data-label="Text"><p>We mentioned that we can do better tokenization if we have regular expressions</p></div> <div data-bbox="27 384 544 411" data-label="List-Group"><ul style="list-style-type: none">• The module <code>re</code> provides regular expression functionality</div> <div data-bbox="769 655 784 669" data-label="Text"><p>19</p></div>	<div data-bbox="1157 178 1276 212" data-label="Section-Header"><h2><code>re.compile</code></h2></div> <div data-bbox="837 279 1596 327" data-label="Text"><p>The first step is to let python compile the regular expression into an internal format that it can use</p></div> <div data-bbox="837 350 1235 476" data-label="Code-Block"><pre>import re p = re.compile('ab*') # can also pass flags to compile: p = re.compile('ab*', re.IGNORECASE)</pre></div> <div data-bbox="837 497 1511 525" data-label="Text"><p><code>p</code> is now a regular expression object, but we still have to know how to use it</p></div> <div data-bbox="1581 655 1596 669" data-label="Text"><p>20</p></div>
<div data-bbox="289 804 521 837" data-label="Section-Header"><h2>Performing matches</h2></div> <div data-bbox="27 913 693 1054" data-label="List-Group"><ul style="list-style-type: none">• <code>match</code> — determines if the RE matches at the <i>beginning</i> of the string• <code>search</code> — determines if the RE matches anywhere in the string<ul style="list-style-type: none">– <code>group</code> — find the part of the string which matched– <code>start</code>, <code>end</code>, <code>span</code> — find the coverage of the match within the string</div> <div data-bbox="27 1075 423 1251" data-label="Code-Block"><pre>s = "this is my string #123" p = re.compile('[0-9]+') # or '\d+' print p.match(s) # None m = p.search(s) print m.group() # '123' print m.span() # (19, 22)</pre></div> <div data-bbox="769 1283 784 1297" data-label="Text"><p>21</p></div>	<div data-bbox="1105 804 1330 837" data-label="Section-Header"><h2>Finding all matches</h2></div> <div data-bbox="842 921 1318 1001" data-label="List-Group"><ul style="list-style-type: none">• <code>findall</code> — returns a list of matches• <code>finditer</code> — allows you to iterate over all matches</div> <div data-bbox="837 1022 1526 1199" data-label="Code-Block"><pre>s = "45 is a number, not number 2, but this is my string #123" p = re.compile('\d+') p.findall(s) # ['45', '2', '123'] iterator = p.finditer(s) for match in iterator: print match.span() # (0, 2) \n (27, 28) \n (53, 56)</pre></div> <div data-bbox="1581 1283 1596 1297" data-label="Text"><p>22</p></div>
<div data-bbox="131 1434 678 1467" data-label="Section-Header"><h2>Putting it into one command (i.e., no compile)</h2></div> <div data-bbox="27 1535 784 1583" data-label="Text"><p>If you are only going to use the RE once, it might be more readable to put it into one line</p></div> <div data-bbox="27 1617 537 1644" data-label="List-Group"><ul style="list-style-type: none">• It's kind of annoying to use <code>compile</code> and then <code>search</code></div> <div data-bbox="27 1665 336 1738" data-label="Code-Block"><pre>s = "this is my string #123" m = re.search('[0-9]+', s) print m.group()</pre></div> <div data-bbox="769 1908 784 1923" data-label="Text"><p>23</p></div>	<div data-bbox="1125 1434 1308 1467" data-label="Section-Header"><h2>Splitting strings</h2></div> <div data-bbox="837 1535 1248 1562" data-label="Text"><p>We can use REs to tell us how to split up text</p></div> <div data-bbox="842 1591 1205 1671" data-label="List-Group"><ul style="list-style-type: none">• Can specify how many splits to make, too• Here, <code>\W</code> refers to non-word characters</div> <div data-bbox="837 1692 1624 1818" data-label="Code-Block"><pre>>>> p = re.compile('\W+') >>> p.split('This is a test, short and sweet, of split().') ['This', 'is', 'a', 'test', 'short', 'and', 'sweet', 'of', 'split', ''] >>> p.split('This is a test, short and sweet, of split().', 3) ['This', 'is', 'a', 'test, short and sweet, of split().']</pre></div> <div data-bbox="1581 1908 1596 1923" data-label="Text"><p>24</p></div>

Capturing parentheses

We can use capturing parentheses to tell us exactly what the items were that we split on

```
>>> p = re.compile('\W+')
>>> p2 = re.compile('(\W+)')
>>> p.split('This... is a test.')
['This', 'is', 'a', 'test', '']
>>> p2.split('This... is a test.')
['This', '... ', 'is', ' ', 'a', ' ', 'test', ' ', '']
```

25

Search and replace

Often we want to not only find a regular expression, but replace it with something else—we use sub for that

```
>>> p = re.compile(' (blue|white|red)')
>>> p.sub('colour', 'blue socks and red shoes')
'colour socks and colour shoes'
>>> p.sub('colour', 'blue socks and red shoes', count=1)
'colour socks and red shoes'
```

26

Greediness

One very important thing to know about this RE matching is that it is *greedy*: it'll match the longest possible thing that it can

```
>>> s = '<html><head><title>Title</title>'
>>> len(s)
32
>>> print re.match('<.*>', s).span()
(0, 32)
>>> print re.match('<.*>', s).group()
<html><head><title>Title</title>
```

There is a way to do non-greedy matches, but I'll let you read the documentation for that.

27

REs for tokenization

```
def tokenize(line):
    list = line.split()
    tokens = []
    for item in list:
        while re.match('\W', item):
            # non-alphanumeric item at beginning of item
            tokens.append(item[0])
            item = item[1:]
        # to maintain order, we use temp
        temp = []
        while re.search('\W$', item):
            # non-alphanumeric item at end of item
            temp.append(item[-1])
            item = item[:-1]
```

28

```
# Contraction handling
if item == "can't":
    tokens.append("can")
    tokens.append("n't")
# other n't words:
elif re.search("n't", item):
    tokens.append(item[:-3])
    tokens.append(item[-3:])
# other words with apostrophes ('s, 'll, etc.)
elif re.search("'", item):
    wordlist = item.split("'")
    tokens.append(wordlist[0])
    tokens.append("'" + wordlist[1])
# no apostrophe, i.e., normal word:
else:
    tokens.append(item)
tokens.extend(temp[::-1])
return tokens
```

29