

NLP Lab Session Week 2
January 23, 2011

Starting a Python and an NLTK Session

Open a Python 2.7 IDLE (Python GUI) window by going to All Programs->Python 2.7 -> IDLE (Python GUI).

You will probably want to work by having the IDLE window open for testing NLTK and a browser window open with these instructions. You may also want to have a separate tab or window open to the NLTK book: <http://www.nltk.org/book/>, where these examples are taken from Chapter 1.

In the following, examples for you to try are given following the Python Idle prompt of >>>. You can copy and paste the Python example into the Idle window, or you can type the example in.

Python and NLTK Resources

Python WikiBooks: http://en.wikibooks.org/wiki/Python_Programming
Python tutorial from python.org: <http://docs.python.org/2/tutorial/>

NLTK book: <http://nltk.org/book/>

NLTK API: <http://nltk.googlecode.com/svn/trunk/doc/api/frames.html>

Counting Frequencies of Words

To get started, we'll repeat some of the steps that we did in the last lab to obtain some text examples from the Gutenberg Corpus available in the NLTK.

```
>>> import nltk
```

For purposes of this lab, again we will work with the first book, Jane Austen's "Emma". First, we get the fileid for Emma and then get the raw text as a string.

```
>>> file1 = nltk.corpus.gutenberg.fileids() [0]
>>> emmatext = nltk.corpus.gutenberg.raw(file1)
>>> len(emmatext)
```

Since this is quite long, we can view part of it, e.g. the first 120 characters

```
>>> emmatext[:120]
```

NLTK has several tokenizers available to break the raw text into tokens; we will use one that separates by white space and also by special characters (punctuation):

```
>>> emmatokens = nltk.wordpunct_tokenize(emmatext)
>>> len(emmatokens)
```

```
>>> emmatokens[:50]
```

Let's decide that we want to count upper case words such as "They" to be the same as a lower case "they", and to simplify things, we won't worry about proper names. So we convert all the tokens to their lower case version and call them words:

```
>>> emmawords = [w.lower() for w in emmatokens]
>>> emmawords[:50]
>>> len(emmawords)
```

[Side note on Python: Here we defined emmawords by using a Python list comprehension. It would be equivalent to define emmawords by starting with an empty list and using a for loop that kept appending the lower case of each word in emmatokens .

Note the special syntax of Python for a multi-line statement: The first line must be followed by an extra ":", and each succeeding line must be indented by some number of spaces (as long as they are all indented by the **same** number of spaces.) Just hit enter to finish the statement.

```
>>> emmawords = []
>>> for w in emmatokens:
    emmawords.append(w.lower())
```

```
>>> emmawords[:50]
End note.]
```

Python Frequency Dictionary

We will start to create a frequency list using a Python dictionary. These are described in the NLTK book, at the end of Chapter 5. We will create a dictionary that has (key, value) pairs where the key is the word and the value is the frequency. We start with an empty dictionary.

```
>>> emmadict = {}
```

We will then write a for loop to go through the words and count them. To do that, if the word is already in the dictionary, we add 1 to its count and if not, we start the count at 1. Note that we indent for the "for" loop and then we indent some more for the body of the "if" and the "else" parts.

```
>>> for w in emmawords:
    if w in emmadict:
        emmadict[w] += 1
    else:
        emmadict[w] = 1
```

How many unique words were counted? We can get the number of keys in the dictionary.

```
>>> len(emmadict.keys())
7344
```

We can look up individual words and print the frequency:

```
>>> emmadict['the']
5201
>>> emmadict['of']
4291
>>> emmadict['a']
3129
```

We don't really want to print the frequencies of all 7,000 words, but we can print 30 of them, using the first 30 keys, but we note that the keys of a dictionary don't occur in any particular order:

```
>>> for word in emmadict.keys()[:30]:
    print word, emmadict[word]
```

Now what we really want is the list of word frequencies in ranked order, i.e. the most frequent first. We can do that in Python, but it's a little awkward, so NLTK has conveniently defined a class to do that for us.

NLTK Frequency Distributions

NLTK has a set of functions that use a data structure called a Frequency Distribution, `FreqDist`.

This structure is an extension of the Python dictionary structures. We can import it from the `nltk.probability` module.

```
>>> from nltk.probability import FreqDist
```

This class allows you to make a Frequency Distribution just by initializing it with a list of words. It will do all the counting for you and create a distribution in which the set of keys are all the words, and the set of values are the frequency (count) of each word. The `keys()` function produces the list of words in order of decreasing frequency.

```
>>> fdist = FreqDist(emmawords)
>>> fdist.keys()[:50]
```

We can treat the frequency distribution just like a Python dictionary and we can look at the frequencies of individual words:

```
>>> fdist['emma']
>>> fdist['the']
```

Or we can use a for loop to look at the frequencies of the first words, and this time the keys are sorted in the order of decreasing values.

```
>>> for word in fdist.keys()[:30]:  
    print word, fdist[word]
```

Try it out:

The text that we first imported from NLTK book are already separated into lists of words. Create a frequency distribution for the words in text1 from the NLTK book (Moby Dick) by applying FreqDist directly to text1. For example:

```
>>> from nltk.book import *  
>>> mbdist = FreqDist(text1)
```

Use text1 as shown or pick one of the other texts and make a FreqDist. Print out some portion of the keys and pick several words and look at the frequencies.

Exercise to submit this week:

Go to Blackboard and find the Discussion for the second week exercises. Create a post in which you:

State which text you used in the “try it out”, give some of the top keywords and give the frequencies of two of the words.