

Introduction to Natural Language Processing with NLTK

Steven Bird Edward Loper Ewan Klein

University of Melbourne, AUSTRALIA

University of Pennsylvania, USA

University of Edinburgh, UK

minor edits by Anoop Sarkar – any mistakes are his fault

Python: Key Features

- ▶ simple yet powerful, shallow learning curve
- ▶ object-oriented: encapsulation, re-use
- ▶ scripting language, facilitates interactive exploration
- ▶ excellent functionality for processing linguistic data
- ▶ extensive standard library, incl graphics, web, numerical processing
- ▶ downloaded for free from <http://www.python.org/>

Python Example

```
import sys
for line in sys.stdin.readlines():
    for word in line.split():
        if word.endswith('ing'):
            print word
```

1. whitespace: nesting lines of code; scope
2. object-oriented: attributes, methods (e.g. line)
3. readable

Comparison with Perl

```
while (<>) {  
    foreach my $word (split) {  
        if ($word =~ /ing$/) {  
            print "$word\n";  
        }  
    }  
}
```

1. syntax is obscure: *what are: <> \$ my split ?*
2. “it is quite easy in Perl to write programs that simply look like raving gibberish, even to experienced Perl programmers”
(Hammond *Perl Programming for Linguists* 2003:47)
3. large programs difficult to maintain, reuse

What NLTK adds to Python

NLTK defines a basic infrastructure that can be used to build NLP programs in Python. It provides:

- ▶ Basic classes for representing data relevant to natural language processing
- ▶ Standard interfaces for performing tasks, such as tokenization, tagging, and parsing
- ▶ Standard implementations for each task, which can be combined to solve complex problems
- ▶ Extensive documentation, including tutorials and reference documentation
- ▶ Large collection of useful language data sets that can be used for non-trivial NLP tasks

Installing Python and NLTK

1. Install Python 2.5.x, Numpy (and optionally, Matplotlib)
2. Install NLTK 0.9, NLTK-Corpora 0.9
3. Set environment variable NLTK_DATA and PATH

For detailed instructions, see:

- ▶ <http://nltk.sourceforge.net/install.html>
- ▶ CDROM: /webpage/install.html