

Chapter 1

Appendix: Projects with NLTK

1.1 Introduction

This document describes a variety of possible natural language processing projects that can be undertaken using NLTK.

The NLTK team welcomes contributions of good student projects, and some past projects (e.g. the Brill and HMM taggers) have been incorporated into the toolkit.

1.2 Project Topics

1.2.1 Computationally Oriented

1. Port the Snowball/TextIndexNG stemmers to NLTK.
2. Implement the TnT statistical tagger in NLTK. <http://www.aclweb.org/anthology/A00-1031>
3. Develop a maximum-entropy POS tagger for NLTK (e.g. see MXPOST)
4. Develop a sentence boundary detector (e.g. see MXTerminator, <http://acl.ldc.upenn.edu/A/A00/A00-1012.pdf> and LingPipe: <http://www.alias-i.com/lingpipe/web/demo-sentence.html>)
5. Develop a chunker that uses transformation-based learning, adapting NLTK's Brill Tagger to chunk tags (see [\[Ramshaw & Marcus, 1995\]](#)).
6. Develop a lexical-chain based WSD system, using the similarity measures defined on WordNet, and evaluate it using the SEMCOR corpus (corpus reader provided in NLTK).
7. Re-implement any NLTK functionality for a language other than English (tokenizer, tagger, chunker, parser, etc). You will probably need to collect suitable corpora, and develop corpus readers.
8. Implement a dependency parser (cf <http://www.seas.upenn.edu/~strctlrn/MSTParser/MSTParser.html>).

9. Create a database of named entities, categorized as: person, location, organization, cardinal, duration, measure, date. Train a named-entity tagger using the NIST IEER data (included with NLTK) and use it to tag more text and collect an expanded set of named entities.
10. Implement a chat-bot that incorporates a more sophisticated dialogue model than `nltk.chat.eliza`.
11. Implement a categorial grammar parser, including semantic representations (see `nltk.contrib.lambek`)
12. Develop a prepositional phrase attachment classifier, using the `ppattach` corpus for training and testing.
13. Develop a program for unsupervised learning of phonological rules, using the method described by Goldwater and Johnson: <http://acl.ldc.upenn.edu/acl2004/sigphon/pdf/goldwater.pdf>
14. Use WordNet to infer lexical semantic relationships on the entries of a Shoebox lexicon for some arbitrary language.
15. Develop a temporal expression identifier (i.e., a system capable of identifying expressions such as *last Christmas, a fortnight ago*), build a temporal expression grounder that will assign specific timestamps to these expressions, e.g. [Day: 25; Month:12; Year: 2005]. Test the accuracy of your system on the TIMEX dataset.
16. Taking the VerbOcean data which captures semantic relationships between verbs (<http://semantics.isi.edu/ocean/verboccean.unrefined.2004-05-20.txt.gz>), generate a semantic network of verb relationships and implement a tree traversal algorithm that can calculate the similarity between two verbs, e.g. *fly* and *crash*. You can find a demo of this system at: <http://falcon.isi.edu/cgi-bin/graph-analysis/view-graph.pl>
17. News stories from different sources often contain contradictory information regarding a particular event such as the number of people killed in an earthquake. Build a numerical expression recognizer and resolver that can identify equality and contradiction between numerical expression such as: *5 adults != 3 children and 2 adults*, but *5 people = 3 children and 2 adults*.
18. The [Gsearch](#) system allows corpora to be searched using a combination of conditions involving lexical content, part-of-speech tag and constituent label; e.g., `gsearch bnc "<tag=NN.* & word=show>" pp` searches for occurrences of the noun *show* followed by a PP. Build a search tool which offers a similar interface. Although Gsearch takes unparsed corpora as input, it would be acceptable to limit your search tool to corpora which have been parsed.

1.2.2 Linguistically Oriented

1. Develop a morphological analyzer for a language of your choice.

2. Develop a coreference resolution system, cf LingPipe <http://www.alias-i.com/lingpipe/web/demo-coref.html>, or an anaphora resolution system, cf MARS <http://clg.wlv.ac.uk/MARS/index.php>
3. Write a soundex function that is appropriate for a language you are interested in. If the language has clusters (consonants or vowels), consider how reliably people can discriminate the second and subsequent member of a cluster. If these are highly confusable, ignore them in the signature. If the *order* of segments in a cluster leads to confusion, normalize this in the signature (e.g. sort each cluster alphabetically, so that a word like `treatments` would be normalized to `rtaemtenst`, before the code is computed). (NB. See `field.html` for more details.)
4. Develop a text classification system which efficiently classifies documents in two or three closely related languages. Consider the discriminating features between languages despite their apparent similarity. Implementation should be evaluated using unseen data.
5. Explore the phonotactic system of a language you are interested in. Compare your findings to a published phonological or grammatical description of the same language.
6. Implement a structured text rendering module which takes linguistic data from a source such as Shoebox and generates XML based lexicon or interlinear text based on user preferences for field exports.
7. Develop a grammatical paradigm generation function which takes some form of tagged text as input and generates paradigm representations of related linguistic features.
8. Build character n-gram models for different languages using the UDHR corpus, and use these to generate hypothetical proper names in these languages (cf. <http://pywordgen.sourceforge.net/>)

1.2.3 Other Sources of Ideas for NLTK Projects

- <http://gate.ac.uk/gate/doc/plugins.html>
- <http://www.alias-i.com/lingpipe/>
- <http://opennlp.sourceforge.net/projects.html>

1.3 Assessment

This section describes the project assessment requirements for *433-460 Human Language Technology* at the University of Melbourne. Project assessment has three components: an oral presentation (5%), a written report (10%), and an implementation (20%).

1.3.1 Oral Presentation

Students will give a 10-minute oral presentation to the rest of the class in the second-last week of semester. This will be evaluated for the quality of content and presentation:

- presentation (clarity, presentation materials, organization)
- content (defining the task, motivation, data, results, outstanding issues)

1.3.2 Written Report

Students should submit a ~5-page written report, with approximately one page covering each of the following points:

- introduction (define the task, motivation)
- method (any algorithms, data)
- implementation (description, how to run it)
- results (e.g. show some output and discuss)
- evaluation (your critical discussion of the work)

This should be prepared using the Python `docutils` and `doctest` packages. These are easily learnt, and ideally suited for creating reports with embedded program code, and they have been used for all NLTK-Lite documentation. For a detailed example, see the text source for the NLTK tagging chapter ([text](#), [html](#)).

- **Docutils**: an open-source text processing system for processing plain-text documentation into useful formats, such as HTML or LaTeX. It includes `reStructuredText`, the easy to read, easy to use, what-you-see-is-what-you-get plain-text markup language.
- **Doctest***: a standard Python module that searches for pieces of text that look like interactive Python sessions, and then executes those sessions to verify that they work exactly as shown.

1.3.3 Implementation

Marks will be awarded for the basic implementation and for various kinds of complexity, as described below:

- Basic implementation (10%)
 - we are able to run the system
 - we can easily test the system (interface is usable, output is appropriately detailed and clearly formatted)
 - we can easily work out how the system is implemented (understandable code, inline documentation; you can assume we read the report first)
 - the system implements NLP algorithms (i.e. relevant to the subject, re-using existing NLP algorithms wherever possible instead of reinventing the wheel)
 - the NLP algorithms are correctly implemented
- Complexity (10%)
 - the system implements a non-trivial problem
 - the system combines multiple HLT components as appropriate
 - appropriate training data is used (effort in obtaining and preparing the data will be considered)

- the system permits exploration of the problem domain and the algorithms (e.g. through appropriate parameterization)
- a range of system configurations/modifications are explored (e.g. classifiers trained and tested using different parameters)

About this document...

This chapter is a draft from *Natural Language Processing* [<http://nltk.org/book.html>], by Steven Bird, Ewan Klein and Edward Loper, Copyright © 2008 the authors. It is distributed with the *Natural Language Toolkit* [<http://nltk.org/>], Version 0.9.5, under the terms of the *Creative Commons Attribution-Noncommercial-No Derivative Works 3.0 United States License* [<http://creativecommons.org/licenses/by-nc-nd/3.0/us/>].

This document is