

PyNLPI: Python Natural Language Processing Library

Maarten van Gompel

2-12-2011



First things first

PyNLPL is pronounced as...



Introduction

What is PyNLPI?

Python Natural Language Processing Library

- A collection of custom-made Python modules usable in Natural Language Processing
- Very modular setup
- Reusable object-oriented modules, prevents “reinventing the wheel” for common tasks
- Using PyNLPI enables you to more quickly write NLP tools, as you need not start from scratch.

Installation

Installation

PyNLPI is on github: <http://github.com/proycon/pynlpl>

To obtain it: \$ git clone

<https://github.com/proycon/pynlpl>

Questions and Answers (1/4)

Q: Why reinvent the wheel yourself and not use for example NLTK?

A: Firstly because there are many customised modules not present in NLTK, such as modules for dealing with FoLiA, D-Coi, Timbl, Cornetto, DutchSemCor. Secondly, because reimplementing things myself was a good learning process to better understand certain algorithms.

Questions and Answers (2/4)

Q: How did PyNLPI came to be?

A: Often code is (and should be) modular and reusable in the future. Whenever that is the case, I put it into PyNLPI.

Questions and Answers (3/4)

Q: Where is PyNLPI used?

A: In almost everything I write: PBMBMT, Valkuil, the DutchSemCor Supervised-WSD system heavily rely on PyNLPI.

Questions and Answers (4/4)

Q: Why Python?

A: Elegant, modern and powerful scripting language, short development time. Great for text processing, easy to learn. Substantial user-base and 3rd party libraries available.

Packages and modules in PyNLPI (1/3)

Packages and modules in PyNLPI (1/3)

- `pynlpl.statistics` – Module containing classes and functions for statistics
- `pynlpl.evaluation` – Module for evaluation and experimentation, such as computation of precision/recall, creation of confusion matrices, etc.. Also contains abstract experiment classes and Wrapped Progressive Sampling.
- `pynlpl.datatypes` – Module containing data types
- `pynlpl.search` – Module containing search algorithms
- `pynlpl.textprocessors` – Module containing text processors

Packages and modules in PyNLPI (2/3)

Packages and modules in PyNLPI (2/3)

- `pynlpl.formats` – Contains modules for reading/writing specific file formats
 - `pynlpl.formats.timbl` – Module for reading Timbl output format
 - `pynlpl.formats.sonar` – Module for reading the SoNaR corpus (D-Coi XML)
 - `pynlpl.formats.folia` - Module for reading/manipulating/writing FoLiA XML
 - `pynlpl.formats.giza` - Module containing class for reading GIZA A3 alignment
 - `pynlpl.formats.moses` - Module containing class for reading phrase translation tabel
 - `pynlpl.formats.cgn` – `pynlpl.formats.dutchsemcor`

Packages and modules in PyNLPI (3/3)

- `pynlpl.clients` – Contains network clients for various services.
 - `pynlpl.clients.cornetto` – Client to connect to Cornetto webservice
 - `pynlpl.clients.frogclient` – Client to connect to Frog server
- `pynlpl.lm` – Language Models
 - `pynlpl.lm.lm` – Contains simple language model
 - `pynlpl.lm.srilm` – SRILM module
 - `pynlpl.lm.server` – Generic LM Server

Using pynlpl.statistics: FrequencyList and Distribution

```
1 >>> from pynlpl.statistics import FrequencyList,  
2   Distribution  
3 >>> freqlist = FrequencyList()  
4 >>> freqlist.append(['It', 'is', 'what', 'is', 'is'])  
5 {'It': 1, 'it': 1, 'is': 2, 'what': 1}  
6 >>> freqlist['is']  
7 2  
8 >>> dist = Distribution(freqlist)  
9 >>> dist  
10 {'It': 0.2, 'it': 0.2, 'is': 0.4, 'what': 0.2}  
11 >>> dist['is']  
12 0.4  
13 >>> dist.entropy()  
14 1.9219280948873623
```

Creating N-grams with `pynlpl.textprocessors.Windower`

```
1 >>> from pynlpl.textprocessors import Windower
2 >>> s = ['It', 'is', 'what', 'it', 'is']
3 >>> list(Windower(s,2))
4 [('<begin>', 'It'), ('It', 'is'), ('is', 'what'),
5 ('what', 'it'), ('it', 'is'), ('is', '<end>')]
```

Combining things: Creating a tri-gram frequency list from a FoLiA document

```
1 >>> from pynlpl.formats import folia
2 >>> from pynlpl.statistics import FrequencyList
3 >>> doc = folia.Document(file='/path/to/folia_doc.xml')
4 >>> freqlist = FrequencyList()
5 >>> for trigram in Windower(doc.words(),3, False, False):
6 ...     freqlist.count(trigram)
7 >>> freqlist.save('freqlist.txt')
```

Creating a simple trigram Language Model of a corpus in FoLiA or DCOI XML

```
1 >>> from pynlpl.formats.folia import Corpus
2 >>> from pynlpl.statistics import FrequencyList
3 >>> simplelm = SimpleLanguageModel(3)
4 >>> for doc in Corpus('/path/to/for/example/sonar/'):
5 ...     for sentence in doc.sentences():
6 ...         simplelm.append([word.text() for word in
7 ...             sentence.words()])
7 >>> simplelm.save('sonar.trigram.lm')
```

Using the Frog Client

First start Frog in server mode: `frog --skip=p -S 12345`

```
1 >>> from pynlpl.clients.frogclient import FrogClient
2 >>> client = FrogClient('localhost', 12345)
3 >>> for word, lemma, morph, pos in client.process("Het u
        iswathetis"):
4     ...     print lemma, pos
5     het VM(pers,pron,stan,red,3,ev,onz)
6     zijn WM(pv,tgw,ev)
7     wat VM(vb,pron,stan,vol,3o,ev)
8     het VM(pers,pron,stan,red,3,ev,onz)
9     zijn WM(pv,tgw,ev)
```

Using the evaluation module (1/2)

```
1 >>> from pynlpl.evaluation import ClassEvaluation
2 >>> weather_forecast = ['sun', 'sun', 'rain', 'cloudy']
3 >>> actual_weather = ['cloudy', 'sun', 'rain', 'rain']
4 >>> evl = ClassEvaluation(actual_weather,
5 >>>                 weather_forecast)
6 >>> print evl
7
8
9
10
11
12
13
14
```

		TP	FP	TN	FN	Accuracy	Precision	Recall(TPR)	Specificity (TNR)	F-score
7	<i>sun</i>	1	1	2	0	0.750000	0.500000	1.000000	0.666667	0.666667	
8	<i>cloudy</i>	0	1	2	1	0.500000	0.000000	0.000000	0.666667	0.000000	
9	<i>rain</i>	1	0	2	1	0.750000	1.000000	0.500000	1.000000	0.666667	
10	<i>Accuracy</i>					: 0.5					
11	<i>Recall</i>					(<i>macroav</i>)	: 0.625				
12	<i>Precision</i>					(<i>macroav</i>)	: 0.5				
13	<i>Specificity</i>					(<i>macroav</i>)	: 0.75				
14											

Using the evaluation module (2/2)

```
1 >>> evl.confusionmatrix()
2 {('cloudy', 'sun'): 1, ('rain', 'cloudy'): 1, ('sun', 'sun'):
3     1, ('rain', 'rain'): 1}
4 >> print evl.confusionmatrix()
5 == Confusion Matrix == (hor: goals, vert: observations)
6
7          cloudy   rain   sun
8      cloudy       0     1     0
9      rain         0     1     0
10     sun          1     0     1
```

Complex topics

Search Algorithms

- ➊ Define your search state, a class derived from the abstract class `pynlpl.search.AbstractSearchState`
- ➋ Add methods `expand()` and for informed searches `score()`
- ➌ Instantiate an initial search state
- ➍ Pass this to the search algorithm of your choice, there are several implemented in `pynlpl.search`: `DepthFirstSearch`, `BreadthFirstSearch`, `IterativeDeepening`, `BestFirstSearch`, `BeamSearch`, `HillClimbingSearch`, `StochasticBeamSearch`
- ➎ Obtain the solution(s) and/or path(s)

Experiments

Experiments

You can define your experiment, as a class derived from the abstract class `pynlpl.evaluation.AbstractExperiment`. Overloading methods as `run()`, `start()`

You can then use these with

`pynlpl.evaluation.ExperimentPool` for multi-threaded use and in `pynlpl.evaluation.ParamSearch` and `pynlpl.evaluation.WPSParamSearch` for parameter optimisation.

Conclusion

Contribute!

If you have a modular, re-usable, preferably object-oriented Python module useful for NLP tasks. Consider adding it to PyNLPI!

Conclusion (shameless promotion)

- **Use** PyNLPI if it has modules you can use!
- **Contribute** to PyNLPI with new modules!

