

HW#1 & Getting Started

LING 571 — Deep Processing Techniques for NLP
Shane Steinert-Threlkeld

Department Cluster

- Assignments are **required** to run on department cluster
 - If you don't have a cluster account, request one ASAP!
 - Link to account request form on Canvas or below:
 - <https://cldb.ling.washington.edu/live/accountrequest-form.php>
- You are not required to develop on the cluster, but code must run on it

Department Cluster

- Assignments are **required** to run on department cluster
 - If you don't have a cluster account, request one ASAP!
 - Link to account request form on Canvas or below:
 - <https://cldb.ling.washington.edu/live/accountrequest-form.php>
- You are not required to develop on the cluster, but code must run on it
- ***Reminder: All but most simple tasks must be run via Condor***

Condor

- Parallel computing management system
- All homework will be run via condor
- See [documentation on CLMS wiki](#) for:
 - Construction of condor scripts
 - Link also on course page under “Course Resources”

NLTK

- Most assignments will use NLTK in Python
- **Natural Language ToolKit (NLTK)**
 - Large, integrated, fairly comprehensive
 - Stemmers
 - Taggers
 - Parsers
 - Semantic analysis
 - Corpus samples
 - ...& More
 - Extensively documented
 - Pedagogically Oriented
 - Implementations Strive for Clarity
 - ...sometimes at the expense of efficiency.

NLTK

- nltk.org
 - Online book
 - Demos of software
 - How-Tos for specific components
 - API information, etc.

Python for 571

- We will use Python for this (and all 57x) course
 - Some introductions at: python.org, docs.python.org
 - Orientation tutorial: <https://github.com/shanest/python-tutorial-clms>
- We have provided a *conda virtual environment* for this class on patas
- To invoke on patas / in scripts, just use full path to binary:
 - **`/dropbox/23-24/571/envs/571/bin/python`**
 - See “`/dropbox/23-24/571/python-example`” for an example bash script

Python for 571

- To develop locally:
 - Install Anaconda/miniconda
 - Scp envs/requirements.txt to your machine

```
conda create -n 571 python=3.10
conda activate 571
conda install --file requirements.txt
```

Python & NLTK

- Interactive mode allows experimentation, introspection:

```
patas$ python
```

```
>>> import nltk
```

```
>>> dir(nltk)
```

```
['AbstractLazySequence', 'AffixTagger', 'AlignedSent',  
'Alignment', 'AnnotationTask', 'ApplicationExpression',  
'Assignment', 'BigramAssocMeasures', 'BigramCollocationFinder',  
'BigramTagger', 'BinaryMaxentFeatureEncoding', ...
```

```
>>> help(nltk.AffixTagger)
```

Python & NLTK

- We will make use of some NLTK data resources in this class.
 - If you use the course environment/binary, you will be good to go
- If using NLTK locally, you will need to, from interactive python:

```
>>> import nltk
>>> nltk.download("punkt")
```

Turning In Homework

- Will be using Canvas' file submission mechanism
- Quick how to at:
<https://community.canvaslms.com/docs/DOC-10663-421254353>

Turning In Homework

- Will be using Canvas' file submission mechanism
 - Quick how to at:
<https://community.canvaslms.com/docs/DOC-10663-421254353>
- Homeworks due on **Wednesday** nights

Turning In Homework

- Will be using Canvas' file submission mechanism
 - Quick how to at:
<https://community.canvaslms.com/docs/DOC-10663-421254353>
- Homeworks due on **Wednesday** nights
- 11:59 PM, Pacific Time

Turning In Homework

- Will be using Canvas' file submission mechanism
 - Quick how to at:
<https://community.canvaslms.com/docs/DOC-10663-421254353>
- Homeworks due on **Wednesday** nights
- 11:59 PM, Pacific Time
- Generally, each assignment will include:
 - `readme.{txt|pdf}`
 - `hwX.tar.gz`
 - Where "X" is the assignment number
 - `tar -cvzf hwX.tar.gz <hw_path>`

HW #1

- Read in sentences and corresponding grammar
- Use NLTK to parse those sentences
- Goals:
 - Set up software environment for rest of course
 - Get familiar with NLTK
 - Work with parsers and CFGs

HW #1: Useful Tools

- Loading data:
 - **`nltk.data.load(resource_url)`**
 - Reads in and processes formatted CFG/FCFG/treebank/etc
 - Returns a grammar from CFG
 - **examples:**
 - `nltk.data.load('grammars/sample_grammars/toy.cfg')`
 - `nltk.data.load('file://' + my_grammar_path)`
 - (NB: absolute path!)

HW #1: Useful Tools

- Loading data:
 - **`nltk.data.load(resource_url)`**
 - Reads in and processes formatted CFG/FCFG/treebank/etc
 - Returns a grammar from CFG
 - **examples:**
 - `nltk.data.load('grammars/sample_grammars/toy.cfg')`
 - `nltk.data.load('file://' + my_grammar_path)`
 - (NB: absolute path!)
- Tokenization:
 - **`nltk.word_tokenize(mystring)`**
 - Returns array of tokens in string
 - (This is why you need “punkt”)

HW #1: Useful Tools

- Parsing:
 - `parser = nltk.parse.EarleyChartParser(grammar)`
 - Returns parser based on the grammar
 - `parser.parse(token_list)`
 - Returns iterator of parses:

```
>>> for item in parser.parse(tokens):  
>>>     print(item)
```

```
(S (NP (Det the) (N dog)) (VP (V chased) (NP (Det the) (N cat))))
```