

# Natural Language Processing:

## Assignment 6: Hooking up words and phrases and clauses

Jordan Boyd-Graber

Out: **16. September 2014**

Due: **17. October 2014**

### Introduction

As always, check out the Github repository with the course homework templates:

[git://github.com/ezubacic/cl1-hw.git](https://github.com/ezubacic/cl1-hw.git)

The code for this homework is in the `hw6` directory.

### Problem 1: Treebank Probabilities (5 points)

To warm up, you'll first compute PCFG probabilities from data. You'll create a class that can read in sentences (through the `add_sentence` function and then answer queries using the `query` function.

Add the required functionality to the file `treebank.py`. You can see how the code will be called in Listing [1](#).

### Problem 2: Shift-Reduce Parsers (15 points)

Next, you'll create sequences of shift-reduce actions that can produce dependency parses from a string of words. You'll implement the `transition_sequence` function in the `oracle.py` file. Your code will be much simpler if you use [generators](#).

An example of the shift-reduce sequence is given in Listing [2](#).

### Problem 3: Eisner Algorithm (40 points)

The goal of this section is for you to build code that can, given a parsing score function and a sentence, discovers the best untyped dependency tree that explains that sentence.

Listing 1: Estimating production probabilities from Treebank data

```
1 >>> tb_probs = PcfgEstimator()
2 >>> from nltk.corpus import treebank
3 >>> for ii in treebank.parsed_sents():
4 ...   tb_probs.add_sentence(ii)
5 ...
6 >>> tb_probs.query('NN', 'man')
7 0.0009114385538508279
```

You'll also need some data files that are too big to put in Github. Please download them here:

- `tb_counts.words`
- `tb_counts.tag`

We're using a very simple score function built from the `Dependencies over Time` corpus. It does not use information about the direction or the part of speech of the corpus; it only uses the words involved (which is the head and which is the dependent). You can assume that all the values are log probabilities (i.e., less than or equal to zero).

You are responsible for completing the following functions in the `EisnerParser` class:

- **`initialize_chart`**: Create a chart that you can fill in later as you parse a sentence. Start with the base case of dynamic programming, giving single spans a score of 0.0.
- **`fill_chart`**: After your chart has been initialized, this function should fill in all the cells in the chart with their appropriate values (this is the hardest part!).
- **`reconstruct`**: Given the coordinates of a span, return an iterator over all of the edges that are part of the best parse of that span. This will require you to store backpointers in your `fill_chart` function.

## Extra Credit

Extra credit is available by providing a custom score function with the function `custom_sf` that performs better than the score function I provided. You

Listing 2: Creating a shift-reduce sequence from a dependency parse

```
1 >>> from nltk.corpus import treebank
2 >>> print(" ".join(conll2007.sents('esp.train')[1]))
3 El Banco_Central_Europeo - BCE - fijo el cambio oficial
4     del euro en los 0,9355_dolares .
5 >>> s = conll2007.parsed_sents('esp.train')[1]
6 >>> for x in transition_sequence(s):
7 ...     print(x.pretty_print(s))
8 ...
9 S
10 S
11 l (Banco_Central_Europeo, El)
12 S
13 S
14 l (BCE, -)
15 S
16 r (BCE, -)
17 r (Banco_Central_Europeo, BCE)
18 S
19 l (fijo, Banco_Central_Europeo)
20 S
21 S
22 l (cambio, el)
23 S
24 r (cambio, oficial)
25 S
26 S
27 r (del, euro)
28 r (cambio, del)
29 r (fijo, cambio)
30 S
31 S
32 S
33 l (0,9355_dlares, los)
34 r (en, 0,9355_dolares)
35 r (fijo, en)
36 S
37 r (fijo, .)
38 r (None, fijo)
39 S
```

can store information (up to 1MB) in the optional upload `sf.dat`, which you can use in the `custom_sf` function).

- You will get no points if you do not describe how the score function is computed
- You will get at most 3 points if you use the Treebank data
- You will get at most 10 points if you use data other than the Treebank data

## Turning in Your Assignment

As usual, submit your completed code to [Moodle](#). Make sure:

- You do not change filenames or function names
- You do not change the API of the functions
- Make sure you pass the included unit tests
- Do not add `print` statements to the code you turn in
- Add your name as a comment to all files you turn in