

## COMS W4705 - Natural Language Processing - Homework 2

**Due: 5:40pm on Wednesday, March 6th**

**Total Points: 100**

You are welcome to discuss the problems with other students but you must turn in your own work. Please review the academic honesty policy for this course (at the end of the syllabus page).

Create a .zip or .tgz archive containing any of the files you submit. Upload that single file to Courseworks.

The file you submit should use the following naming convention:

YOURUNI\_homework2.[zip | tgz]. For example, my uni is db2711, so my file should be named db2711\_homework2.zip or db2711\_homework2.tgz.

As a reminder, any assignments submitted late will incur a 20 point penalty. No submissions will be accepted later than 4 days after the submission deadline.

### Analytical Component (40 pts)

Write up your solution in a single .pdf document. Image files and Microsoft Word documents will not be accepted. If you must upload a scan or photo converted into .pdf, make sure that the file size does not exceed 1MB. Name your file written.txt or written.pdf and include it in the zip file you upload to Courseworks.

#### Problem 1 (10 pts) - PCFGs and HMMs

Both PCFGs and HMMs can be seen as generative models that produce a sequence of POS tags and words with some probability (of course the PCFG will generate even more structure, but it will *also* generate POS tags and words).

(a) Consider the grammar specified in Problem 2 below, and the sentence "*they are baking potatoes*". For each sequence of POS tags that is possible for this sentence according to the grammar, what is the *joint* probability  $P(\text{tags}, \text{words})$  according to the PCFG? Hint: consider all parses for the sentence -- you may want to work on problem 2 first.

(b) Design an HMM that produces the same joint probability  $P(\text{tags}, \text{words})$  as the PCFG for each of the possible tag sequences for the sentence in part (a). Note: Your HMM does **not** have to assign 0 probabilities to tag sequences that are not possible according to the PCFG.

(c) In general, is it possible to translate *any* PCFG into an HMM that produces the identical *joint* probability  $P(\text{tags}, \text{words})$  as the PCFG (i.e. not just for a single sentence)? Explain how or why not. No formal proof is necessary. Hint: This has nothing to do with probabilities, but with language classes.

## Problem 2 (10 pts) - Earley Parser

Consider the following probabilistic context free grammar.

$S \rightarrow NP VP$	[1.0]
$NP \rightarrow Adj NP$	[0.3]
$NP \rightarrow PRP$	[0.1]
$NP \rightarrow N$	[0.6]
$VP \rightarrow V NP$	[0.8]
$VP \rightarrow Aux V NP$	[0.2]
$PRP \rightarrow they$	[1.0]
$N \rightarrow potatoes$	[1.0]
$Adj \rightarrow baking$	[1.0]
$V \rightarrow baking$	[0.5]
$V \rightarrow are$	[0.5]
$Aux \rightarrow are$	[1.0]

(a) Using this grammar, show how the **Earley algorithm** would parse the following sentence.

*they are baking potatoes*

Write down the complete parse chart. The chart should contain  $n+1$  entries where  $n$  is the length of the sentence. Each entry  $i$  should contain *all* parser items generated by the parser that end in position  $i$ .

You can ignore the probabilities for part (a).

(b) Write down *all* parse trees for the sentence and grammar from problem 2 and compute their probabilities according to the PCFG.

## Problem 3 (10 pts) - CKY parsing

(a) Convert the grammar from problem 2 into an equivalent grammar in Chomsky Normal Form (CNF). Write down the new grammar. Also explain what the general rule is for dealing with

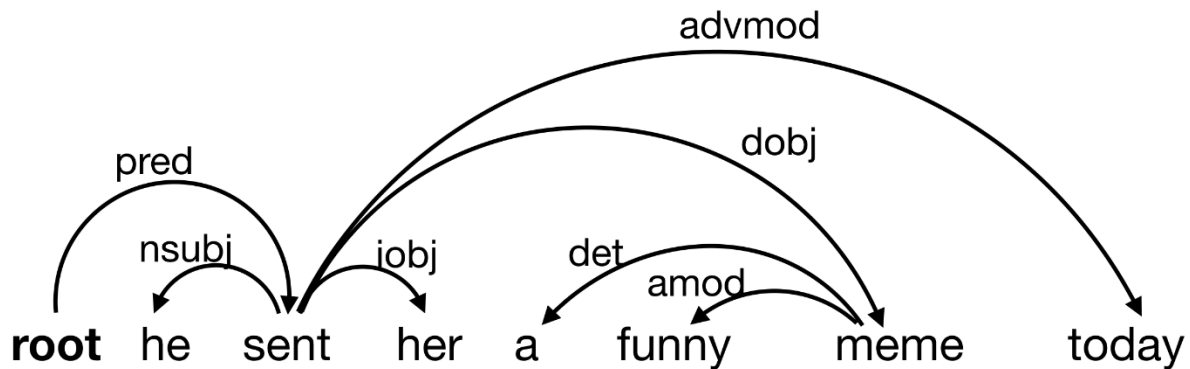
1. Rules of the form  $A \rightarrow B$  (i.e. a single nonterminal on the right hand side).
2. Rules with three or more nonterminals on the right hand side (e.g.  $A \rightarrow B C D E$ ).

You do not have to deal with the case in which terminals and non-terminals are mixed in a rule right-hand side. You also do not have to convert the probabilities. Hint: Think about adding new nonterminal symbols.

(b) Using your grammar, fill the CKY parse chart as shown in class and show all parse trees.

#### Problem 4 (10 pts) - Transition Based Dependency Parsing

Consider the following dependency graph.



Write down the sequence of transitions that an arc-standard dependency parser would have to take to generate this dependency tree from the initial state

$([\text{root}]_\sigma, [\text{he}, \text{sent}, \text{her}, \text{a}, \text{funny}, \text{meme}, \text{today}]_\beta, \{\}_A)$

Also write down the state resulting from each transition.

## Programming Component - Parsing with Context Free Grammar (60 pts)

The instructions below are fairly specific and it is okay to deviate from implementation details. **However: You will be graded based on the functionality of each function. Make sure the function signatures (function names, parameter and return types/data structures) match exactly the description in this assignment.**

Please make sure you are developing and running your code using Python 3.

### Introduction

In this assignment you will implement the **CKY algorithm** for CFG and PCFG parsing. You will also practice retrieving parse trees from a parse chart and working with such tree data structures.

*The files for this project are inside **hw2.zip** which should be found in the same folder as this pdf.*

### Python files:

To get you started, there are three Python files for this project.

1. cky.py will contain your parser and currently contains only scaffolding code.
2. grammar.py contains the class Pcfg which represents a PCFG grammar (explained below) read in from a grammar file.

3. `evaluate_parser.py` contains a script that evaluates your parser against a test set.

### Data files:

You will work with an existing PCFG grammar and a small test corpus.

The main data for this project has been extracted from the ATIS (Air Travel Information Services) subsection of the Penn Treebank. ATIS is originally a spoken language corpus containing user queries about air travel. These queries have been transcribed into text and annotated with Penn-Treebank phrase structure syntax.

The data set contains sentences such as "*what is the price of flights from indianapolis to memphis .*"

There were 576 sentences in total, out of which 518 were used for training (extracting the grammar and probabilities) and 58 for test. The data set is obviously tiny compared to the entire Penn Treebank and typically that would not be enough training data. However, because the domain is so restricted, the extracted grammar is actually able to generalize reasonably well to the test data.

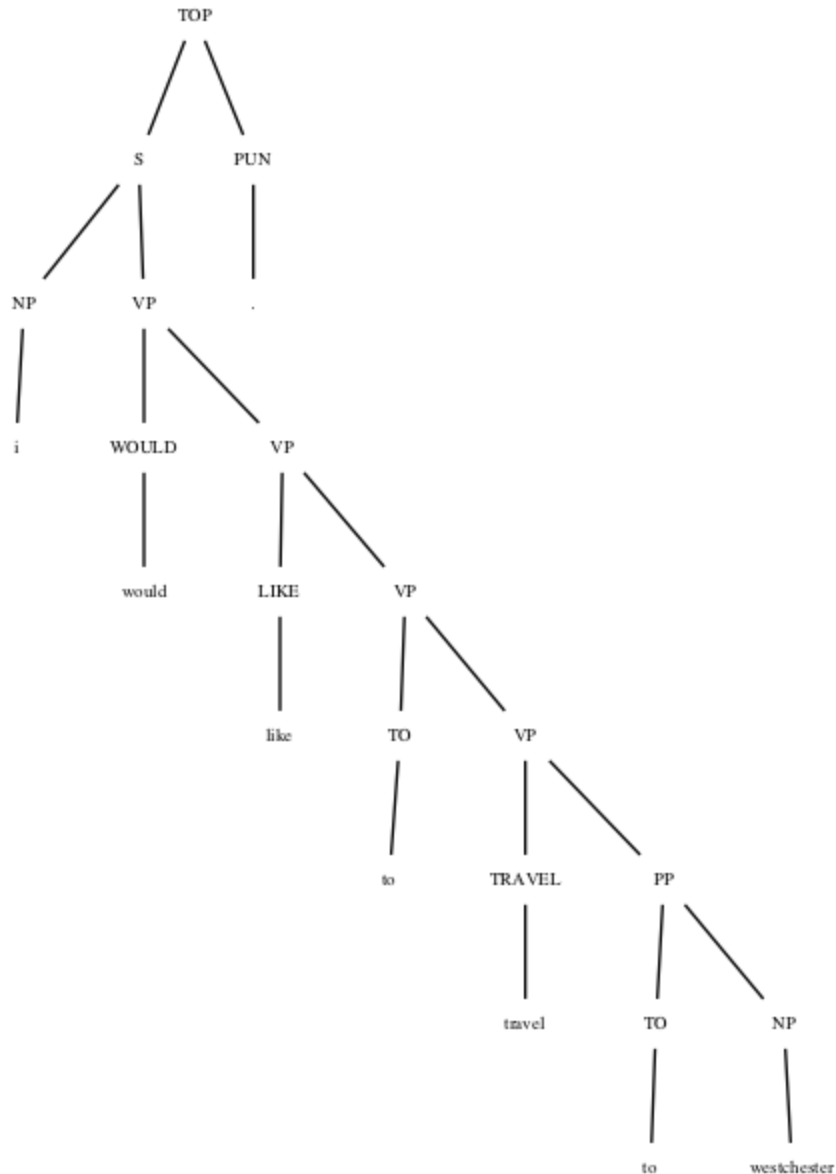
There are 2 data files:

`atis3.pcfg` - contains the PCFG grammar (980 rules)

`atis3_test.ptb` - contains the test corpus (58 sentences).

Take a look at these files and make sure you understand the format. The tree structures in `atis3_test.ptb` are a little different from what we have seen in class. Consider the following example from the test file:

```
(TOP (S (NP i) (VP (WOULD would) (VP (LIKE like) (VP (TO to) (VP (TRAVEL travel) (PP (TO to) (NP westchester))))))) (PUN .))
```



Note that there are no part of speech tags. In some cases phrases like NP directly project to the terminal symbol (NP -> westchester). In other cases, nonterminals for a specific word were added (TRAVEL -> travel).

The start-symbol for the grammar (and therefore the root for all trees) is "TOP". This is the result of an automatic conversion to make the tree structure compatible with the grammar in Chomsky Normal Form.

While you are working on your parser, you might find it helpful to additionally create a small toy grammar (for example, the one in the lecture slides) that you can try on some hand written test cases, so that you can verify by hand that the output is correct.

## Part 1 - reading the grammar and getting started

Take a look at `grammar.py`. The class *Pcfg* represents a PCFG grammar in chomsky normal form. To instantiate a *Pcfg* object, you need to pass a file object to the constructor, which contains the data. For example:

```
with open('atis3.pcfg','r') as grammar_file:
    grammar = Pcfg(grammar_file)
```

You can then access the instance variables of the *Pcfg* instance to get information about the rules:

```
>>> grammar.startsymbol
'TOP'
```

The dictionary *lhs\_to\_rules* maps left-hand-side (lhs) symbols to lists of rules. For example, we will want to look up all rules of the form `PP -> ??`

```
>>> grammar.lhs_to_rules['PP']
[('PP', ('ABOUT', 'NP'), 0.00133511348465), ('PP', ('ADVP', 'PPBAR'), 0.00133511348465), ('PP', ('AFTER', 'NP'), 0.0253671562083), ('PP', ('AROUND', 'NP'), 0.00667556742323), ('PP', ('AS', 'ADJP'), 0.00133511348465), ... ]
```

Each rule in the list is represented as (lhs, rhs, probability) triple. So the first rule in the list would be

*PP -> ABOUT NP* with PCFG probability 0.00133511348465.

The *rhs\_to\_rules* dictionary contains the same rules as values, but indexed by right-hand-side. For example:

```
>>> grammar.rhs_to_rules[('ABOUT', 'NP')]
[('PP', ('ABOUT', 'NP'), 0.00133511348465)]

>>> grammar.rhs_to_rules[('NP', 'VP')]
[('NP', ('NP', 'VP'), 0.00602409638554), ('S', ('NP', 'VP'), 0.694915254237), ('SBAR', ('NP', 'VP'), 0.166666666667), ('SQBAR', ('NP', 'VP'), 0.289156626506)]
```

### TODO:

- Write the method *verify\_grammar*, that checks that the grammar is a valid PCFG in CNF. You need to check that the rules have the right format (note all nonterminal symbols are upper-case) and that all probabilities for the same lhs symbol sum to 1.0. Hint: For improved numeric accuracy, use `math.fsum` to compute the sum.
- Then change the main section of `grammar.py` to read in the grammar, print out a confirmation if the grammar is a valid PCFG in CNF or print an error message if it is not. You should now be able to run `grammar.py` on grammars and verify that they are well formed for the CKY parser.

## Part 2 - Membership checking with CKY

The file `cky.py` already contains a class `CkyParser`. When a `CkyParser` instance is created a `grammar` instance is passed to the constructor. The instance variable `grammar` can then be used to access this `Pcfg` object.

**TODO:** Write the method `is_in_language(self, tokens)` by implementing the CKY algorithm. Your method should read in a list of tokens and return `True` if the grammar can parse this sentence and `False` otherwise. For example:

```
>>> parser = CkyParser(grammar)
>>> toks = ['flights', 'from', 'miami', 'to', 'cleveland', '.'] // Or: toks= 'flights
from miami to cleveland '.split()
>>> parser.is_in_language(toks)
True
>>> toks = ['miami', 'flights', 'cleveland', 'from', 'to', '.']
>>> parser.is_in_language(toks)
False
```

While parsing, you will need to access the dictionary `self.grammar.rhs_to_rules`. You can use any data structure you want to represent the parse table (or read ahead to Part 3 of this assignment, where a specific data structure is prescribed).

The ATIS grammar actually overgenerates a lot, so many unintuitive sentences can be parsed. You might want to create a small test grammar and test cases. Also make sure that this method works for grammar with different start symbols.

## Part 3 - Parsing with backpointers

The parsing method in part 2 can identify if a string is in the language of the grammar, but it does not produce a parse tree. It also does not take probabilities into account. You will now extend the parser so that it retrieves the most probable parse for the input sentence, given the PCFG probabilities in the grammar. The lecture slides on parsing with PCFG will be helpful for this step.

**TODO:** Write the method `parse_with_backpointers(self, tokens)`. You should modify your CKY implementation from part 2, but use (and return) specific data structures. The method should take a list of tokens as input and returns a) the parse table b) a probability table. Both objects should be constructed during parsing. They replace whatever table data structure you used in part 2.

The two data structures are somewhat complex. They will make sense once you understand their purpose.

The first object is **parse table containing backpointers**, represented as a dictionary (this is more convenient in Python than a 2D array). The keys of the dictionary are spans, for example `table[(0,3)]` retrieves the entry for span 0 to 3 from the chart. The values of the dictionary should be dictionaries that map nonterminal symbols to backpointers. For example: `table[(0,3)]['NP']` returns the backpointers to the table entries that were used to create the NP phrase over the span 0 and 3. For example, the value

of `table[(0,3)]['NP']` could be `((("NP",0,2),("FLIGHTS",2,3)))`. This means that the parser has recognized an NP covering the span 0 to 3, consisting of another NP from 0 to 2 and FLIGHTS from 2 to 3. The split recorded in the table at `table[(0,3)]['NP']` is the one that results in the ***most probable parse for the span [0,3] that is rooted in NP***. Terminal symbols in the table could just be represented as strings. For example the table entry for `table[(2,3)]['FLIGHTS']` should be "flights".

The **second object is similar, but records log probabilities** instead of backpointers. For example the value of `probs[(0,3)]['NP']` might be -12.1324. This value represents the log probability of the *best parse tree (according to the grammar)* for the span 0,3 that results in an NP.

Your `parse_with_backpointers(self, tokens)` method should be called like this:

```
>>> parser = CkyParser(grammar)
>>> toks = ['flights', 'from', 'miami', 'to', 'cleveland', '.']
>>> table, probs = parser.parse_with_backpointers(toks)
```

During parsing, when you fill an entry on the backpointer parse table and iterate through the possible splits for a (span/nonterminal) combination, that entry on the table will contain the back-pointers for the the current-best split you have found so far. For each new possible split, you need to check if that split would produce a higher log probability. If so, you update the entry in the backpointer table, as well as the entry in the probability table.

After parsing has finished, the table entry `table[0,len(toks)][grammar.startsymbol]` will contain the best backpointers for the left and right subtree under the root node. `probs[0,len(toks)][grammar.startsymbol]` will contain the total log-probability for the best parse.

`cky.py` contains two test functions `check_table_format(table)` and `check_prob_format(probs)` that you can use to make sure the two table data structures are formatted correctly. Both functions should return True. Note that passing this test does not guarantee that the content of the tables is correct, just that the data structures are probably formatted correctly.

```
>>> check_table_format(table)
True
>>> check_prob_format(probs)
True
```

## Part 4 - Retrieving a parse tree

You now have a working parser, but in order to evaluate its performance we still need to reconstruct a parse tree from the backpointer table returned by `parse_with_backpointers`.

Write the function `get_tree(chart, i, j, nt)` which should return the parse-tree rooted in non-terminal `nt` and covering span `i,j`.



For example

```
>>> parser = CkyParser(grammar)
>>> toks = ['flights', 'from', 'miami', 'to', 'cleveland', '.']
>>> table, probs = parser.parse_with_backpointers(toks)
>>> get_tree(table, 0, len(toks), grammar.startsymbol)
('TOP', ('NP', ('NP', 'flights'), ('NPBAR', ('PP', ('FROM', 'from'), ('NP', 'miami'))
, ('PP', ('TO', 'to'), ('NP', 'cleveland')))), ('PUN', '.'))
```

Note that the intended format is the same as the data in the treebank. Each tree is represented as tuple where the first element is the parent node and the remaining elements are children. Each child is either a tree or a terminal string.

Hint: Recursively traverse the parse chart to assemble this tree.

## Part 5 - Evaluating the Parser

The program `evaluate_parser.py` evaluates your parser by comparing the output trees to the trees in the test data set.

You can run the program like this:

```
python evaluate_parser.py atis3.pcfg atis3_test.ptb
```

Even though this program has been written for you, it might be useful to take a look at how it works.

The program imports your CKY parser class. It then reads in the test file line-by-line. For each test tree, it extracts the yield of the tree (i.e. the list of leaf nodes). It then feeds this list as input to your parser, obtains a parse chart, and then retrieves the predicted tree by calling your `get_tree` method.

It then compares the predicted tree against the target tree in the following way (this is a standard approach to evaluating constituency parsers, called PARSEVAL):

First it obtains a set of the spans in each tree (including the nonterminal label). It then computes precision, recall, and F-score (which is the harmonic mean between precision and recall) between these two sets. The script finally reports the coverage (percentage of test sentences that had *any* parse), the average F-score for parsed sentences, and the average F-score for all sentences (including 0 f-scores for unparsed sentences).

My parser implementation produces the following result on the atis3 test corpus:

Coverage: 67%, Average F-score (parsed sentences): 0.95, Average F-score (all sentences): 0.64

### **What you need to submit**

cky.py  
evaluate\_parser.py  
grammar.py  
written.pdf

Do not submit the data files.

Pack these files together in a zip or tgz file as described on top of this page. **Please follow the submission instructions precisely!**