

CS236299 Project Segment 3: Parsing – The CKY Algorithm

June 27, 2023

```
[1]: # Please do not change this cell because some hidden tests might depend on it.
import os

# Otter grader does not handle ! commands well, so we define and use our
# own function to execute shell commands.
def shell(commands, warn=True):
    """Executes the string `commands` as a sequence of shell commands.

    Prints the result to stdout and returns the exit status.
    Provides a printed warning on non-zero exit status unless `warn`
    flag is unset.
    """
    file = os.popen(commands)
    print (file.read().rstrip('\n'))
    exit_status = file.close()
    if warn and exit_status != None:
        print(f"Completed with errors. Exit status: {exit_status}\n")
    return exit_status

shell("""
ls requirements.txt >/dev/null 2>&1
if [ ! $? = 0 ]; then
    rm -rf .tmp
    git clone https://github.com/cs236299-2023-spring/project3.git .tmp
    mv .tmp/requirements.txt ./
    rm -rf .tmp
fi
pip install -q -r requirements.txt
""")
```

```
[2]: # Initialize Otter
import otter
grader = otter.Notebook()
```

1 236299 - Introduction to Natural Language Processing

1.1 Project 3: Parsing – The CKY Algorithm

Constituency parsing is the recovery of a labeled hierarchical structure, a *parse tree* for a sentence of a natural language. It is a core intermediary task in natural-language processing, as the meanings of sentences are related to their structure.

In this project, you will implement the CKY algorithm for parsing strings relative to context-free grammars (CFG). You will implement versions for both non-probabilistic context-free grammars (CFG) and probabilistic grammars (PCFG) and apply them to the parsing of ATIS queries.

The project is structured into five parts:

1. Finish a CFG for the ATIS dataset.
2. Implement the CKY algorithm for *recognizing* grammatical sentences, that is, determining whether a parse exists for a given sentence.
3. Extend the CKY algorithm for *parsing* sentences, that is, constructing the parse trees for a given sentence.
4. Construct a probabilistic context-free grammar (PCFG) based on a CFG.
5. Extend the CKY algorithm to PCFGs, allowing the construction of the most probable parse tree for a sentence according to a PCFG.

2 Setup

```
[3]: # Download needed files and scripts
import wget
os.makedirs('data', exist_ok=True)
os.makedirs('scripts', exist_ok=True)
# ATIS queries
wget.download("https://raw.githubusercontent.com/nlp-236299/data/master/ATIS/
↳train.nl", out="data/")
# Corresponding parse trees
wget.download("https://raw.githubusercontent.com/nlp-236299/data/master/ATIS/
↳train.trees", out="data/")
wget.download("https://raw.githubusercontent.com/nlp-236299/data/master/ATIS/
↳test.trees", out="data/")

# Code for comparing and evaluating parse trees
wget.download("https://raw.githubusercontent.com/nlp-236299/data/master/scripts/
↳trees/evalb.py", out="scripts/")
wget.download("https://raw.githubusercontent.com/nlp-236299/data/master/scripts/
↳trees/transform.py", out="scripts/")
wget.download("https://raw.githubusercontent.com/nlp-236299/data/master/scripts/
↳trees/tree.py", out="scripts/")
```

```
[3]: 'scripts//tree (1).py'
```

```
[4]: import shutil

import nltk

import sys

from collections import defaultdict, Counter

from nltk import treetransforms
from nltk.grammar import ProbabilisticProduction, PCFG, Nonterminal
from nltk.tree import Tree

from tqdm import tqdm

# Import functions for transforming augmented grammars
sys.path.insert(1, './scripts')
import transform as xform

[5]: ## Debug flag used below for turning on and off some useful tracing
DEBUG = False
```

3 A custom ATIS grammar

To parse, we need a grammar. In this project, you will use a hand-crafted grammar for a fragment of the ATIS dataset. The grammar is written in a “semantic grammar” style, in which the nonterminals tend to correspond to semantic classes of phrases, rather than syntactic classes. By using this style, we can more closely tune the grammar to the application, though we lose generality and transferability to other applications. The grammar will be used again in the next project segment for a question-answering application.

We download the grammar to make it available.

```
[6]: if not os.path.exists('./data/grammar_distrib3'):
    wget.download(
        "https://raw.githubusercontent.com/nlp-236299/data/master/ATIS/
        ↪grammar_distrib3", out="data/")
if os.path.exists('./data/grammar_distrib3') and (not os.path.exists('./data/
    ↪grammar')):
    shutil.copy('./data/grammar_distrib3', './data/grammar')
```

Take a look at the file `data/grammar_distrib3` that you’ve just downloaded. The grammar is written in a format that extends the NLTK format expected by `CFG.fromstring`. We’ve provided functions to make use of this format in the file `scripts/transform.py`. You should familiarize yourself with this format by checking out the documentation in that file.

We made a copy of this grammar for you as `data/grammar`. This is the file you’ll be modifying in the next section. You can leave it alone for now.

As described there, we can read the grammar in and convert it into NLTK's grammar format using the provided `xform.read_augmented_grammar` function.

```
[7]: atis_grammar_distrib, _ = xform.read_augmented_grammar(
      "grammar_distrib3", path="data")
```

To verify that the ATIS grammar that we distributed is working, we can parse a sentence using a built-in NLTK parser. We'll use a tokenizer built with NLTK's tokenizing apparatus.

```
[8]: # Tokenizer
tokenizer = nltk.tokenize.RegexpTokenizer('\d+|[\w-]+|\$[\d\.]+|\S+')

def tokenize(string):
    return tokenizer.tokenize(string.lower())

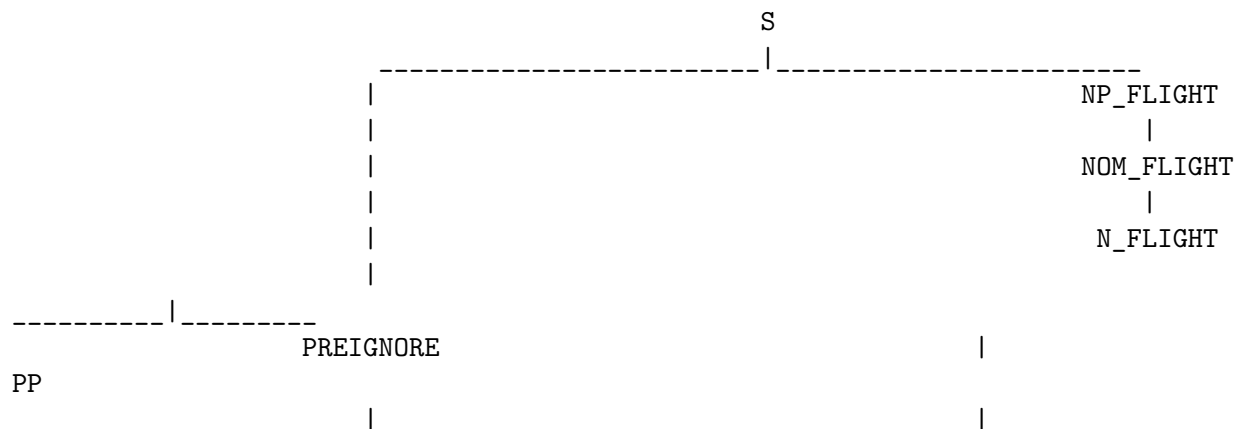
# Demonstrating the tokenizer
# Note especially the handling of `11pm` and hyphenated words.
print(tokenize("Are there any first-class flights at 11pm for less than $3.50?
↪"))
```

```
['are', 'there', 'any', 'first-class', 'flights', 'at', '11', 'pm', 'for',
'less', 'than', '$3.50', '?']
```

```
[9]: # Test sentence
test_sentence_1 = tokenize("show me the flights before noon")

# Construct parser from distribution grammar
atis_parser_distrib = nltk.parse.BottomUpChartParser(atis_grammar_distrib)

# Parse and print the parses
parses = atis_parser_distrib.parse(test_sentence_1)
for parse in parses:
    parse.pretty_print()
```



3.1 Testing the coverage of the grammar

We can get a sense of how well the grammar covers the ATIS query language by measuring the proportion of queries in the training set that are parsable by the grammar. We define a `coverage` function to carry out this evaluation.

Warning: It may take a long time to parse all of the sentence in the training corpus, on the order of 30 minutes. You may want to start with just the first few sentences in the corpus. The `coverage` function below makes it easy to do so, and in the code below we just test coverage on the first 50 sentences.

```
[10]: ## Read in the training corpus
with open('data/train.nl') as file:
    training_corpus = [tokenize(line) for line in file]

[11]: def coverage(recognizer, corpus, n=0):
    """Returns the proportion of the first `n` sentences in the `corpus`
    that are recognized by the `recognizer`, which should return a boolean.
    `n` is taken to be the whole corpus if n is not provided or is
    non-positive.
    """
    n = len(corpus) if n <= 0 else n
    parsed = 0
    total = 0
    for sent in tqdm(corpus[:n]):
        total += 1
        try:
            parses = recognizer(sent)
        except:
            parses = None
        if parses:
            parsed += 1
        elif DEBUG:
            print(f"failed: {' '.join(sent)}")
    if DEBUG: print(f"{parsed} of {total}")
    return parsed/total

[12]: coverage(lambda sent: 0 < len(list(atis_parser_distrib.parse(sent))), # trick
    ↪ for turning parser into recognizer
    training_corpus, n=50)
```

```
100%|          | 50/50 [00:00<00:00, 574.32it/s]
```

```
[12]: 0.0
```

Sadly, you'll find that the coverage of the grammar is extraordinarily poor. That's because it is missing crucial parts of the grammar, especially phrases about *places*, which play a role in essentially every ATIS query. You'll need to complete the grammar before it can be useful.

3.2 Part 1: Finish the CFG for the ATIS dataset

Consider the following query:

```
[13]: test_sentence_2 = tokenize("show me the united flights from boston")
```

You'll notice that the grammar we distributed doesn't handle this query because it doesn't have a subgrammar for airline information ("united") or for places ("from boston").

```
[14]: len(list(atis_parser_distrib.parse(test_sentence_2)))
```

```
[14]: 0
```

Follow the instructions in the grammar file `data/grammar` to add further coverage to the grammar. (You can and should leave the `data/grammar_distrib3` copy alone and use it for reference.)

We'll define a parser based on your modified grammar, so we can compare it against the distributed grammar. Once you've modified the grammar, this test sentence should have at least one parse.

You can search for "TODO" in `data/grammar` to find the two places to add grammar rules.

```
[15]: atis_grammar_expanded, _ = xform.read_augmented_grammar("grammar", path="data")
      atis_parser_expanded = nltk.parse.BottomUpChartParser(atis_grammar_expanded)

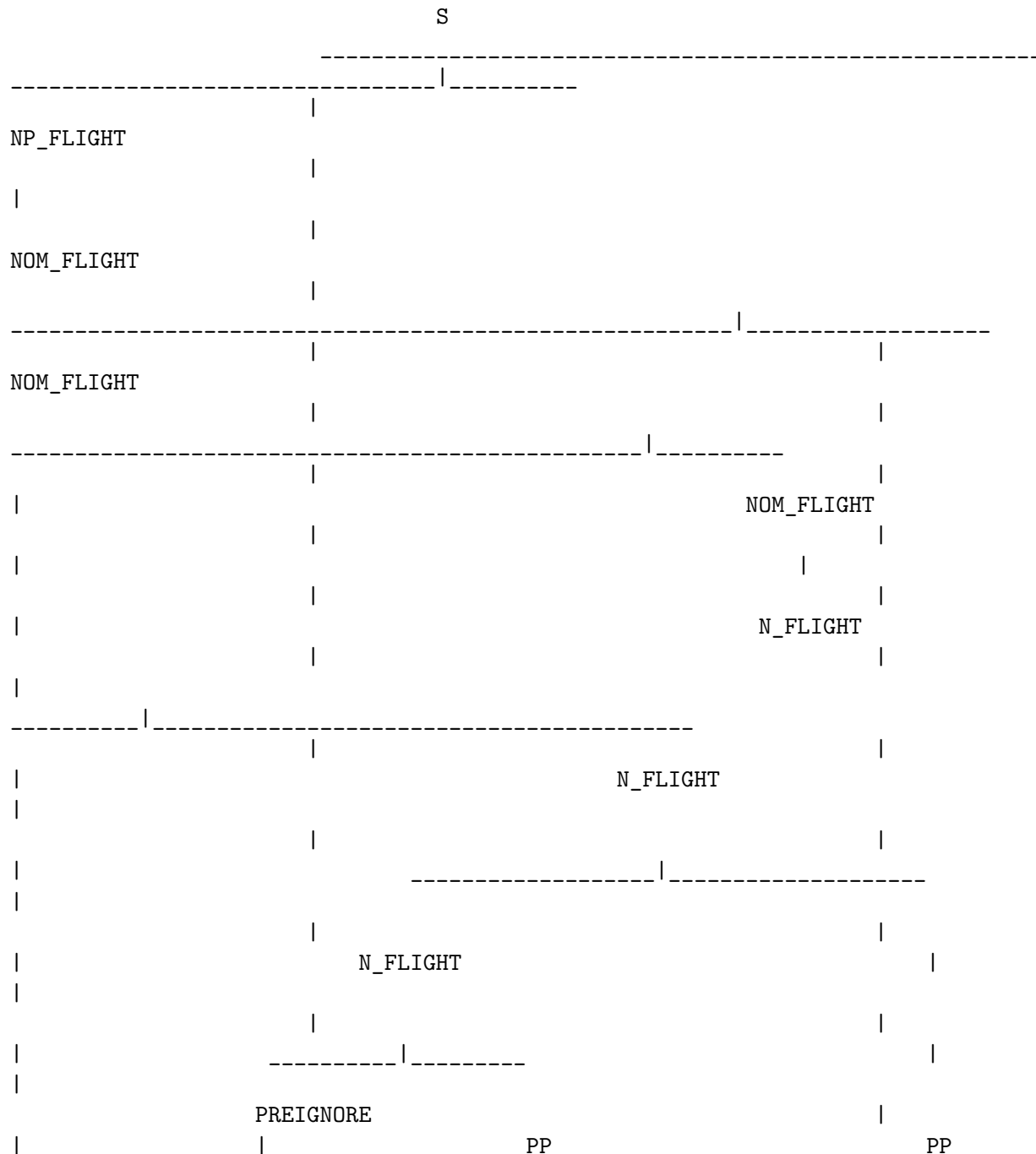
      parses = [p for p in atis_parser_expanded.parse(test_sentence_2)]
      for parse in parses:
          parse.pretty_print()
```

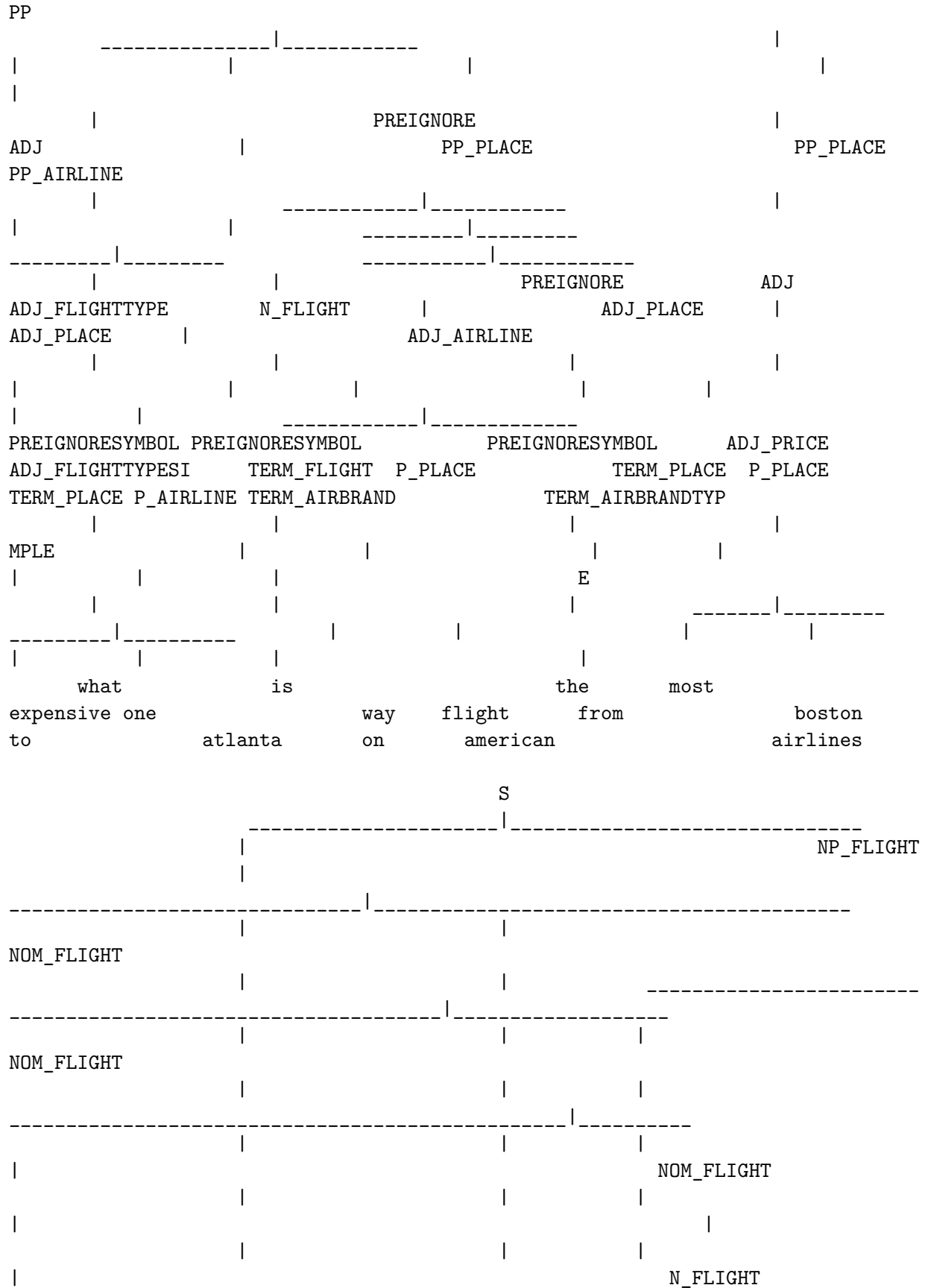
```

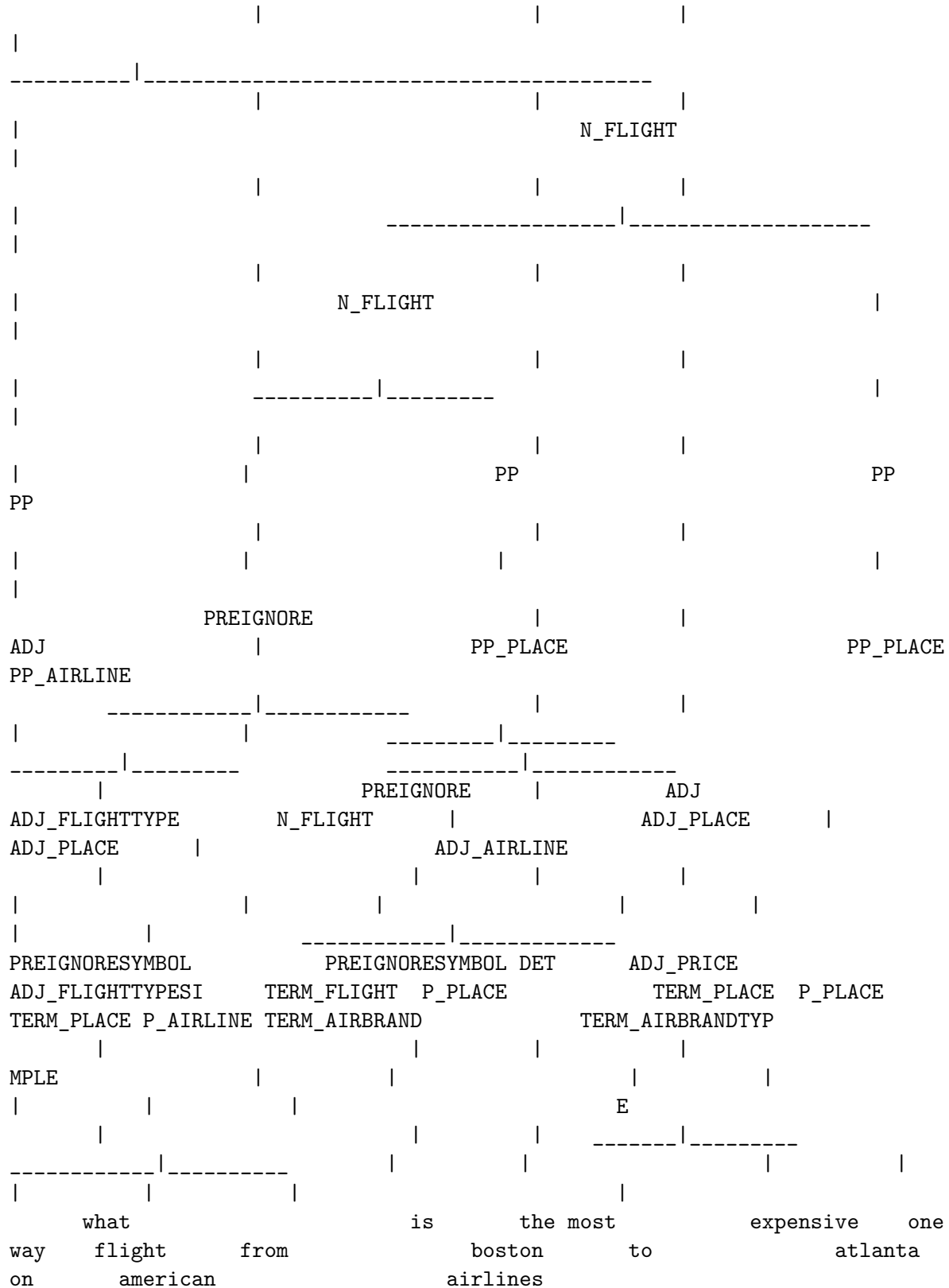
                                                     S
-----|-----
      |
NP_FLIGHT
      |
|
      |
NOM_FLIGHT
      |
-----|-----
      |
NOM_FLIGHT
      |
|
      |
N_FLIGHT
      |
-----|-----
      PREIGNORE
PP
      |
-----|-----
```


show me the united flights
from boston

```
[16]: test_sentence_3 = tokenize("what is the most expensive one way flight from_
    ↪ boston to atlanta on american airlines")
parses = [p for p in atis_parser_expanded.parse(test_sentence_3)]
for parse in parses:
    parse.pretty_print()
```







Once you're done adding to the grammar, to check your grammar, we'll compute the grammar's

coverage of the ATIS training corpus as before. **This grammar should be expected to cover about half of the sentences in the first 50 sentences, and a third of the entire training corpus.**

```
[17]: coverage(lambda sent: 0 < len(list(atis_parser_expanded.parse(sent))), # trick ↵  
    ↪for turning parser into recognizer  
    training_corpus, n=50)
```

```
100%|      | 50/50 [00:00<00:00, 460.51it/s]
```

```
[17]: 0.5
```

4 CFG recognition via the CKY algorithm

Now we turn to implementing recognizers and parsers using the CKY algorithm. We start with a recognizer, which should return `True` or `False` if a grammar does or does not admit a sentence as grammatical.

4.1 Converting the grammar to CNF for use by the CKY algorithm

The CKY algorithm requires the grammar to be in Chomsky normal form (CNF). That is, only rules of the forms

$$A \rightarrow BC$$

$$A \rightarrow a$$

are allowed, where A, B, C are nonterminals and a is a terminal symbol.

However, in some downstream applications (such as the next project segment) we want to use grammar rules of more general forms, such as $A \rightarrow BCD$. Indeed, the ATIS grammar you've been working on makes use of the additional expressivity beyond CNF.

To satisfy both of these constraints, we will convert the grammar to CNF, parse using CKY, and then convert the returned parse trees back to the form of the original grammar. We provide some useful functions for performing these transformations in the file `scripts/transform.py`, already loaded above and imported as `xform`.

To convert a grammar to CNF:

```
cnf_grammar, cnf_grammar_wunaries = xform.get_cnf_grammar(grammar)
```

To convert a tree output from CKY back to the original form of the grammar:

```
xform.un_cnf(tree, cnf_grammar_wunaries)
```

We pass into `un_cnf` a version of the grammar before removing unary nonterminal productions, `cnf_grammar_wunaries`. The `cnf_grammar_wunaries` is returned as the second part of the returned value of `get_cnf_grammar` for just this purpose.

```
[18]: atis_grammar_cnf, atis_grammar_wunaries = xform.  
    ↪get_cnf_grammar(atis_grammar_expanded)  
    assert(atis_grammar_cnf.is_chomsky_normal_form())
```

In the next sections, you'll write your own recognizers and parsers based on the CKY algorithm that can operate on this grammar.

4.2 Part 2: Implement a CKY recognizer

Implement a *recognizer* using the CKY algorithm to determine if a sentence `tokens` is parsable. The labs and J&M Chapter 13, both of which provide appropriate pseudo-code for CKY, should be useful references here.

Hint: Recall that you can get the production rules of a grammar using `grammar.productions()`.

Throughout this project segment, you should use `grammar.start()` to get the special start symbol from the grammar instead of using `S`, since some grammar uses a different start symbol, such as `TOP`.

```
[19]: from typing import List, Tuple

class CKY():
    def __init__(self, grammar: nltk.CFG):
        """
        Initialize the CKY parser with the given grammar.

        Args:
            grammar (nltk.CFG): Context-Free Grammar object.

        """
        self.grammar: nltk.CFG = grammar
        self._table = None
        self._backptrs = None
        self._words = None
        self._N = 0
        self._default_entry: type = set
        self._backptr_entry: type = lambda: defaultdict(set)

    def _init_table(self, tokens: List[str]):
        """
        Initialize the parsing table and other necessary variables.

        Args:
            tokens (List[str]): List of input tokens.

        """
        words = [''] + tokens
        self._N = len(words) - 1
        self._table = defaultdict(self._default_entry)
        self._backptrs = defaultdict(self._backptr_entry)
        self._words = words
```

```

def get_table(self):
    """
    Get the parsing table.

    Returns:
        List[List[set]]: Parsing table.

    """
    assert self._table is not None
    return [[self._table[i, j] if j > i else '---' for j in range(self._N + 1)]
    ↪1)] for i in range(self._N + 1)]

def get_backpointers(self):
    """
    Get the backpointers.

    Returns:
        List[List[defaultdict(set)]]: Backpointers.

    """
    assert self._backptrs is not None
    return [[self._backptrs[i, j] if j > i else '---' for j in range(self.
    ↪_N + 1)] for i in range(self._N + 1)]

def parse(self, tokens: List[str]) -> Tuple[nltk.Tree, bool]:
    """
    Parse the given list of tokens and return the parse tree and whether
    ↪the sentence is recognized.

    Args:
        tokens (List[str]): List of input tokens.

    Returns:
        Tuple[nltk.Tree, bool]: Parse tree and recognition status.

    """
    self._init_table(tokens)

    for j in range(1, self._N + 1):
        # Handle the rules of the form A -> w
        self._handle_potential_constituents_for_terminal(self._words[j], j)

        # Handle rules of the form A -> B C
        for length in range(2, j + 1):
            i = j - length
            for split in range(i + 1, j):

```

```

        for B in self._iterate_nonterminals_in_cell(self._table[i,
↪split]):
            for C in self._iterate_nonterminals_in_cell(self.
↪_table[split, j]):
                self.
↪_handle_potential_constituents_for_nonterminals(
                    [B, C], i, split, j)

    # if DEBUG:
    #     print('Final table:\n', matrix(self.get_table()))
    #     print('Backpointers:\n', matrix(self.get_backpointers()))

    return self._get_parse_result()

def _handle_potential_constituents_for_terminal(self, terminal: str, j:
↪int):
    """
    Handle the productions of the form A -> w that match the terminal at
↪[j-1, j].

    Args:
        terminal (str): Terminal symbol.
        j (int): End string position.

    """
    prods: List[nltk.grammar.Production] = self.grammar.productions()

    for product in prods:
        if terminal in product.rhs():
            self._update_table_for_terminal_production(product, j)

    def _handle_potential_constituents_for_nonterminals(self, nonterminals:
↪List[str], i: int, split: int, j: int):
        """
        Handle the productions of the form A -> B C such that B is in table[i,
↪split] and C is in table[split, j].

    Args:
        nonterminals (List[str]): List of nonterminal symbols.
        i (int): Start string position.
        split (int): Split point.
        j (int): End string position.

    """
    assert len(

```

```

        nonterminals) == 2, "Grammar is assumed to be in CNF therefore_
↳there can only be two non-terminals."
        prods: List[nltk.grammar.Production] = self.grammar.productions()

        B, C = nonterminals[0], nonterminals[1]

        for product in prods:
            if len(product.rhs()) != 2:
                # This is not a product of the form A -> B C
                continue

            if B == product.rhs()[0] and C == product.rhs()[1]:
                self._update_table_for_nonterminal_production(
                    product, i, split, j)

    def _build_trees(self) -> Tree:
        """
        Build the parse tree.

        Returns:
            nltk.Tree: Parse tree.

        """
        if self.grammar.start() in self._table[0, self._N]:
            return [Tree.fromstring(tree_str) for tree_str in self.
↳_construct_trees(0, self._N, (self.grammar.start()))]
        else:
            return None

    def _construct_trees(self, i: int, j: int, symbol: nltk.Nonterminal) -> str:
        """
        Recursively construct the parse tree.

        Args:
            i (int): Start string position.
            j (int): End string position.
            symbol (str): Symbol at the current position.

        Returns:
            str: String ready to convert into a tree using Tree.fromstring

        """
        if i + 1 == j:
            # This is a terminal
            return [f"({symbol} {self._words[j]})"]

        backpointers = self._backptrs[i, j][symbol]

```



```

trees = []
for split, left, right in backpointers:
    left_trees = self._construct_trees(i, split, left)
    right_trees = self._construct_trees(split, j, right)
    trees.extend([f"({symbol} {left_tree} {right_tree})"
                  for left_tree in left_trees
                  for right_tree in right_trees])

return trees

# Override these method in order to extend this class:
def _update_table_for_terminal_production(self, production: nltk.grammar.
↳Production, j: int):
    """
    Handles the productions of the form  $A \rightarrow w$  that match the terminal at  $\_$ 
    ↳ $[j-1, j]$ .

    Args:
        production (nltk.grammar.Production): Production rule.
        j (int): End string position.

    """
    self._table[j - 1, j].add(production.lhs())

def _update_table_for_nonterminal_production(self, production: nltk.grammar.
↳Production, i: int, split: int, j: int):
    """
    Handles the productions of the form  $A \rightarrow B C$  such that  $B$  is in table  $[i, \_$ 
    ↳ $split]$  and  $C$  is in table  $[split, j]$ .

    Args:
        production (nltk.grammar.Production): Production rule.
        i (int): Start string position.
        split (int): Split point.
        j (int): End string position.

    """
    self._table[i, j].add(production.lhs())
    self._backptrs[i, j][production.lhs()].add(
        (split, production.rhs()[0], production.rhs()[1]))

def _iterate_nonterminals_in_cell(self, cell):
    """
    Iterate over the nonterminals in a cell.

    Args:
        cell: Cell containing nonterminals.

```

```

        Yields:
            Nonterminal: Nonterminal symbol.

        """
        for nonterminal in cell:
            yield nonterminal

    def _get_parse_result(self):
        return self._build_trees(), self.grammar.start() in self._table[0, self.
↪_N]

```

```

[20]: # TODO - Implement a CKY recognizer
def cky_recognize(grammar, tokens):
    """
    Returns True if and only if the list of tokens `tokens` is admitted
    by the `grammar`.

    Implements the CKY algorithm, and therefore assumes `grammar` is in
    Chomsky normal form.
    """
    assert (grammar.is_chomsky_normal_form())

    cky = CKY(grammar)
    _, is_admitted = cky.parse(tokens)
    return is_admitted

```

You can test your recognizer on a few examples, both grammatical and ungrammatical, as below.

```

[22]: test_sentences = ["show me flights from boston",
                        "show me united flights before noon",
                        "are there any twa flights available tomorrow",
                        "show me flights united are there any"]

for sentence in test_sentences:
    recognized = "+" if cky_recognize(atis_grammar_cnf,
                                     tokenize(sentence)) else "-"
    print(f"{recognized:5}{sentence}")

```

```

+   show me flights from boston
+   show me united flights before noon
+   are there any twa flights available tomorrow
-   show me flights united are there any

```

You can also verify that the CKY recognizer verifies the same coverage as the NLTK parser.

```

[23]: coverage(lambda sent: cky_recognize(atis_grammar_cnf, sent),
               training_corpus, n=50)

```

100% | 50/50 [00:05<00:00, 9.44it/s]

[23]: 0.5

4.3 Part 3: Implement a CKY parser

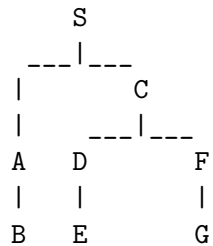
In part 2, you implemented a context-free grammar recognizer. Next, you'll implement a *parser*.

Implement the CKY algorithm for parsing with CFGs as a function `cky_parse`, which takes a grammar and a list of tokens and returns a single parse of the tokens as specified by the grammar, or `None` if there are no parses. You should only need to add a few lines of code to your CKY recognizer to achieve this, to implement the necessary back-pointers. The function should return an NLTK tree, which can be constructed using `Tree.fromstring`.

A tree string will be like this example:

"(S (A B) (C (D E) (F G)))"

which corresponds to the following tree (drawn using `tree.pretty_print()`):



Hint: You may want to extract from a `Nonterminal` its corresponding string. The `Nonterminal.__str__` method or f-string `f'{Nonterminal}'` accomplishes this.

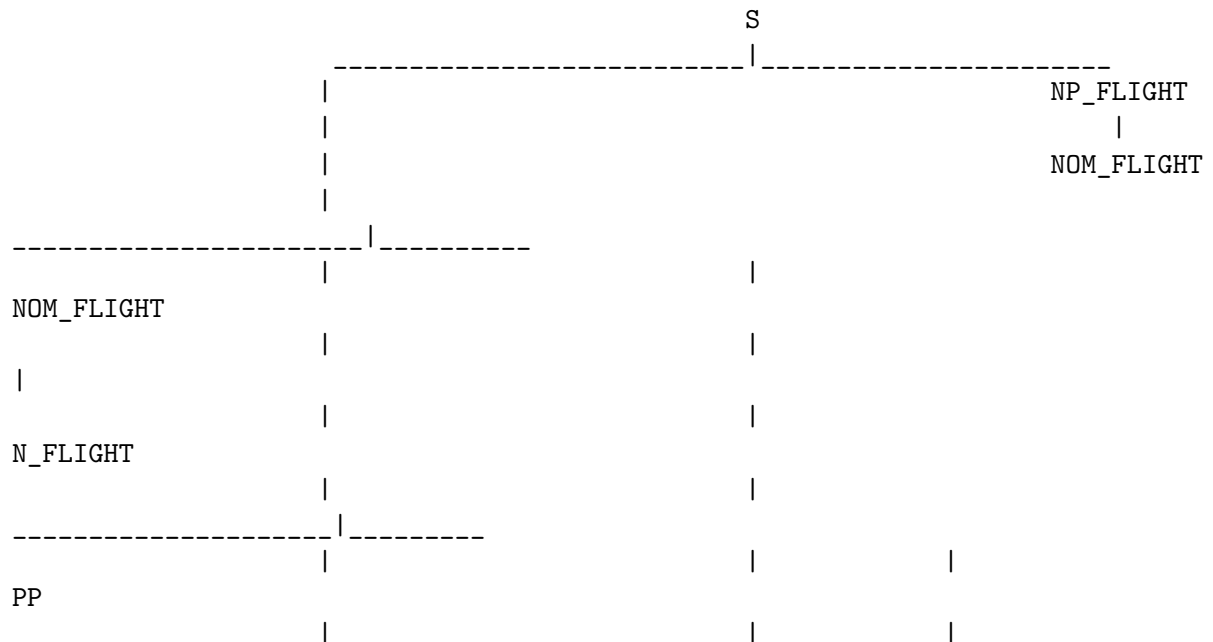
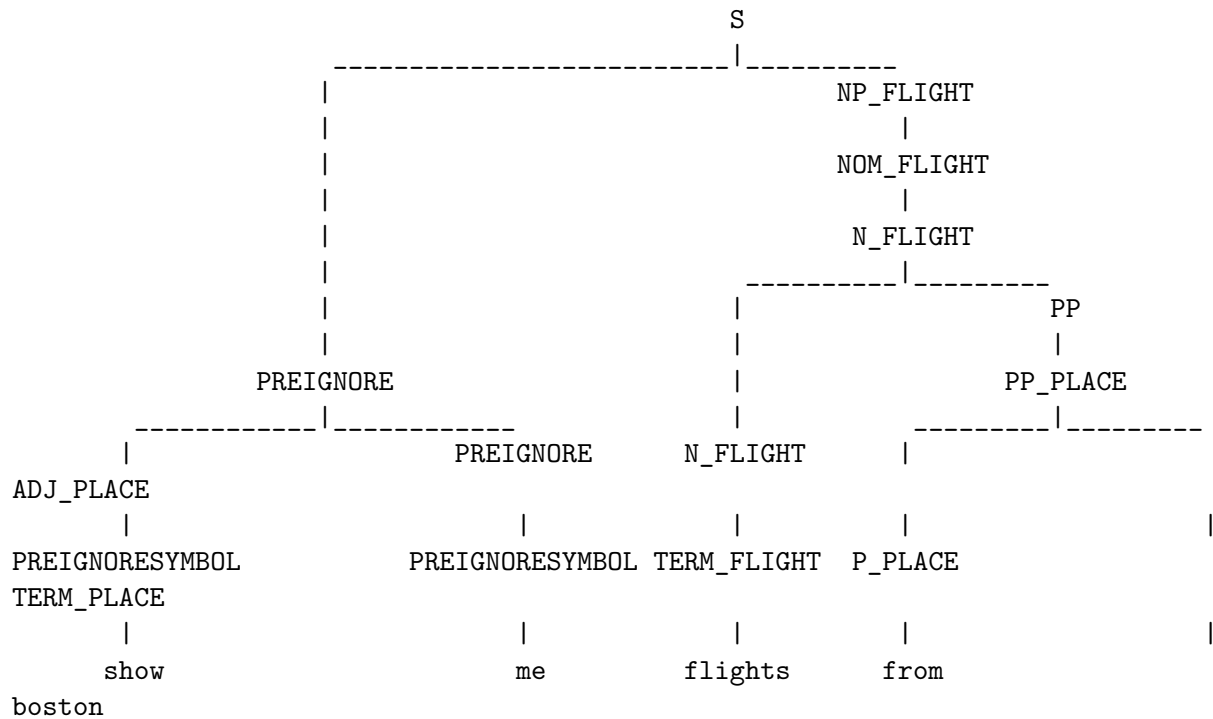
```
[24]: # TODO -- Implement a CKY parser
def cky_parse(grammar, tokens):
    """
    Returns an NLTK parse tree of the list of tokens `tokens` as
    specified by the `grammar`. If there are multiple valid parses,
    return any one of them.

    Returns None if `tokens` is not parsable.
    Implements the CKY algorithm, and therefore assumes `grammar` is in
    Chomsky normal form.
    """
    assert (grammar.is_chomsky_normal_form())

    cky = CKY(grammar)
    trees, _ = cky.parse(tokens)
    return trees[0] if trees else None
```

You can test your code on the test sentences provided above:

```
[25]: for sentence in test_sentences:
      tree = cky_parse(atis_grammar_cnf, tokenize(sentence))
      if not tree:
          print(f"failed to parse: {sentence}")
      else:
          xform.un_cnf(tree, atis_grammar_wunaries)
          tree.pretty_print()
```



		PREIGNORE		ADJ	
PP_TIME					
	----- -----				
----- -----					
		PREIGNORE	ADJ_AIRLINE	N_FLIGHT	
	NP_TIME				
PREIGNORESYMBOL		PREIGNORESYMBOL	TERM_AIRBRAND	TERM_FLIGHT	
P_TIME	TERM_TIME				
show		me	united	flights	
before	noon				
				S	
		----- -----			
				NP_FLIGHT	
				----- -----	
NOM_FLIGHT					
----- -----					
NOM_FLIGHT					
N_FLIGHT					
----- -----					
N_FLIGHT					
----- -----					
	PREIGNORE			ADJ	
PP	PP				
	----- -----				
		PREIGNORE		ADJ_AIRLINE	N_FLIGHT
PP_CLASS	PP_DATE				
PREIGNORESYMBOL		PREIGNORESYMBOL	DET	TERM_AIRBRAND	TERM_FLIGHT
ADJ_CLASS	NP_DATE				

are there any twa flights
available tomorrow

failed to parse: show me flights united are there any

You can also compare against the built-in NLTK parser that we constructed above:

```
[26]: for sentence in test_sentences:
    reparses = [p for p in atis_parser_expanded.parse(tokenize(sentence))]
    predparse = cky_parse(atis_grammar_cnf, tokenize(sentence))
    if predparse:
        xform.un_cnf(predparse, atis_grammar_wunaries)

    print('Reference parses:')
    for reftree in reparses:
        print(reftree)

    print('\nPredicted parse:')
    print(predparse)

    if (not predparse and len(reparses) == 0) or predparse in reparses:
        print("\nSUCCESS!")
    else:
        print("\nOops. No match.")
```

Reference parses:

```
(S
  (PREIGNORE (PREIGNORESYMBOL show) (PREIGNORE (PREIGNORESYMBOL me)))
  (NP_FLIGHT
    (NOM_FLIGHT
      (N_FLIGHT
        (N_FLIGHT (TERM_FLIGHT flights))
        (PP
          (PP_PLACE (P_PLACE from) (ADJ_PLACE (TERM_PLACE boston)))))))))
```

Predicted parse:

```
(S
  (PREIGNORE (PREIGNORESYMBOL show) (PREIGNORE (PREIGNORESYMBOL me)))
  (NP_FLIGHT
    (NOM_FLIGHT
      (N_FLIGHT
        (N_FLIGHT (TERM_FLIGHT flights))
        (PP
          (PP_PLACE (P_PLACE from) (ADJ_PLACE (TERM_PLACE boston)))))))))
```

SUCCESS!

Reference parses:

```
(S
```

```

(PREIGNORE (PREIGNORESYMBOL show) (PREIGNORE (PREIGNORESYMBOL me)))
(NP_FLIGHT
  (NOM_FLIGHT
    (ADJ (ADJ_AIRLINE (TERM_AIRBRAND united)))
    (NOM_FLIGHT
      (N_FLIGHT
        (N_FLIGHT (TERM_FLIGHT flights))
        (PP (PP_TIME (P_TIME before) (NP_TIME (TERM_TIME noon))))))))))

```

Predicted parse:

```

(S
  (PREIGNORE (PREIGNORESYMBOL show) (PREIGNORE (PREIGNORESYMBOL me)))
  (NP_FLIGHT
    (NOM_FLIGHT
      (ADJ (ADJ_AIRLINE (TERM_AIRBRAND united)))
      (NOM_FLIGHT
        (N_FLIGHT
          (N_FLIGHT (TERM_FLIGHT flights))
          (PP (PP_TIME (P_TIME before) (NP_TIME (TERM_TIME noon))))))))))

```

SUCCESS!

Reference parses:

```

(S
  (PREIGNORE
    (PREIGNORESYMBOL are)
    (PREIGNORE (PREIGNORESYMBOL there)))
  (NP_FLIGHT
    (DET any)
    (NOM_FLIGHT
      (ADJ (ADJ_AIRLINE (TERM_AIRBRAND twa)))
      (NOM_FLIGHT
        (N_FLIGHT
          (N_FLIGHT (TERM_FLIGHT flights))
          (PP (PP_CLASS (ADJ_CLASS available))))
          (PP (PP_DATE (NP_DATE tomorrow))))))))

```

Predicted parse:

```

(S
  (PREIGNORE
    (PREIGNORESYMBOL are)
    (PREIGNORE (PREIGNORESYMBOL there)))
  (NP_FLIGHT
    (DET any)
    (NOM_FLIGHT
      (ADJ (ADJ_AIRLINE (TERM_AIRBRAND twa)))
      (NOM_FLIGHT
        (N_FLIGHT

```

```
(N_FLIGHT
  (N_FLIGHT (TERM_FLIGHT flights))
  (PP (PP_CLASS (ADJ_CLASS available))))
(PP (PP_DATE (NP_DATE tomorrow))))))
```

SUCCESS!

Reference parses:

Predicted parse:

None

SUCCESS!

Again, we test the coverage as a way of verifying that your parser works consistently with the recognizer and the NLTK parser.

```
[27]: coverage(lambda sent: cky_parse(atis_grammar_cnf, sent),
              training_corpus, n=50)
```

```
0%|          | 0/50 [00:00<?, ?it/s]100%|      | 50/50 [00:05<00:00,
9.51it/s]
```

```
[27]: 0.5
```

5 Probabilistic CFG parsing via the CKY algorithm

In practice, we want to work with grammars that cover nearly all the language we expect to come across for a given application. This leads to an explosion of rules and a large number of possible parses for any one sentence. To remove ambiguity between the different parses, it's desirable to move to probabilistic context-free grammars (PCFG). In this part of the assignment, you will construct a PCFG from training data, parse using a probabilistic version of CKY, and evaluate the quality of the resulting parses against gold trees.

5.1 Part 4: PCFG construction

Compared to CFGs, PCFGs need to assign probabilities to grammar rules. For this goal, you'll write a function `pcfg_from_trees` that takes a list of strings describing a corpus of trees and returns an NLTK PCFG trained on that set of trees.

We expect you to implement `pcfg_from_trees` directly. You should **not** use the `induce_pcfg` function in implementing your solution.

We want the PCFG to be in CNF format because the probabilistic version of CKY that you'll implement next also requires the grammar to be in CNF. However, the gold trees are not in CNF form, so in this case you will need to convert the gold *trees* to CNF before building the PCFG from them. To accomplish this, you should use the `treetransforms` package from `nltk`, which includes functions for converting to and from CNF. In particular, you'll want to make use of `treetransforms.collapse_unary` followed by `treetransforms.chomsky_normal_form` to convert a tree to its binarized version. You can then get the counts for all of the productions used in the

trees, and then normalize them to probabilities so that the probabilities of all rules with the same left-hand side sum to 1.

We'll use the `pcfg_from_trees` function that you define later for parsing.

To convert an `nlk.Tree` object `t` to CNF, you can use the below code. Note that it's different from the `xform` functions we used before as we are converting *trees*, not *grammars*.

```
treetransforms.collapse_unary(t, collapsePOS=True)
treetransforms.chomsky_normal_form(t) # After this the tree will be in CNF
```

To construct a PCFG with a given start state and set of productions, see `nlk.grammar.PCFG`.

```
[28]: # TODO - Define a function to convert a set of trees to a PCFG in Chomsky
      ↪ normal form.
      # You are not allowed to use any library functions except
      # `treetransforms.collapse_unary` and `treetransforms.chomsky_normal_form`,
      # write the logic by yourself.
def pcfg_from_trees(trees, start=Nonterminal('TOP')):
    """Returns an NLTK PCFG in CNF with rules and counts extracted from a set
    ↪ of trees.

    The `trees` argument is a list of strings in the form interpretable by
    `Tree.fromstring`. The trees are converted to CNF using NLTK's
    `treetransforms.collapse_unary` and `treetransforms.chomsky_normal_form`.

    The `start` argument is the start nonterminal symbol of the returned
    grammar."""
    # Production count: the number of times a given production occurs
    production_counts = {}

    # LHS count: counts the number of times a given LHS occurs
    lhs_counts = {}

    for tree_str in trees:
        tree = Tree.fromstring(tree_str)
        treetransforms.collapse_unary(tree, collapsePOS=True)
        treetransforms.chomsky_normal_form(tree)
        productions = tree.productions()
        for production in productions:
            lhs_counts[production.lhs()] = lhs_counts.get(production.lhs(), 0)
            ↪ + 1
            production_counts[production] = production_counts.get(production,
            ↪ 0) + 1

    probabilistic_productions = [
```

```

        ProbabilisticProduction(prod.lhs(), prod.rhs(),
    ↪prob=production_counts[prod] / lhs_counts[prod.lhs()])
        for prod in production_counts
    ]

    return PCFG(start, probabilistic_productions)

```

We can now train a PCFG on the *train* split `train.trees` that we downloaded in the setup at the start of the notebook.

```

[29]: with open('data/train.trees') as file:
        # Convert the probabilistic productions to an NLTK probabilistic CFG.
        pgrammar = pcfg_from_trees(file.readlines())

        # Verify that the grammar is in CNF
        assert (pgrammar.is_chomsky_normal_form())

```

5.2 Part 5: Probabilistic CKY parsing

Finally, we are ready to implement probabilistic CKY parsing under PCFGs. Adapt the CKY parser from Part 3 to return the most likely parse and its **log probability** (base 2) given a PCFG. Note that to avoid underflows we want to work in the log space. > **Hint:** `production.logprob()` will return the log probability of a production rule `production`.

```

[30]: from typing import Dict

class PCKY(CKY):
    def __init__(self, grammar: PCFG):
        super().__init__(grammar)
        self._default_entry: type = lambda: defaultdict(set)
        self._backptr_entry: type = lambda: defaultdict(set)
        self._prod_to_logprob: Dict[nltk.ProbabilisticProduction, float] = {
            nltk.grammar.Production(prod.lhs(), prod.rhs()): prod.logprob() for ↪
    ↪prod in grammar.productions()
        }

        def _update_table_for_terminal_production(self, production: nltk.grammar.
    ↪ProbabilisticProduction, j: int):
            self._table[j - 1, j][production.lhs()] = production.logprob()

        def _update_table_for_nonterminal_production(self, production: nltk.grammar.
    ↪ProbabilisticProduction, i: int, split: int, j: int):
            assert len(production.rhs()) == 2

            A, (B, C) = production.lhs(), production.rhs()

```

```

        new_prob = production.logprob() + self._table[i, split][B] + self.
↪_table[split, j][C]
        if A not in self._table[i, j] or self._table[i, j][A] < new_prob:
            self._table[i, j][A] = new_prob
            self._backptrs[i, j][A] = [(split, B, C)]

    def _get_parse_result(self):
        tree = self._build_trees()
        logprob = -float('inf')
        if tree is not None:
            logprob = sum([self._prod_to_logprob[prod] for prod in tree[0].
↪productions()])

        return tree[0] if tree else None, logprob

```

```

[31]: # TODO - Implement a CKY parser under PCFGs
def cky_parse_probabilistic(grammar, tokens):
    """Returns the NLTK parse tree of `tokens` with the highest probability
    as specified by the PCFG `grammar` and its log probability as a tuple.

    Returns (None, -float('inf')) if `tokens` is not parsable.
    Implements the CKY algorithm, and therefore assumes `grammar` is in
    Chomsky normal form.
    """

    assert (grammar.is_chomsky_normal_form())

    pcky = PCKY(grammar)
    tree, logprob = pcky.parse(tokens)
    return tree, logprob

```

As an aid in debugging, you may want to start by testing your implementation of probabilistic CKY on a much smaller grammar than the one you trained from the ATIS corpus. Here's a little grammar that you can play with.

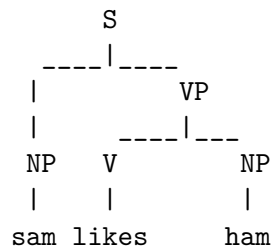
Hint: By “play with”, we mean that you can change the grammar to try out the behavior of your parser on different test grammars, including ambiguous cases.

```

[32]: grammar = PCFG.fromstring("""
    S -> NP VP [1.0]
    VP -> V NP [1.0]
    PP -> P NP [1.0]
    NP -> 'sam' [.3]
    NP -> 'ham' [.7]
    V -> 'likes' [1.0]
    """)

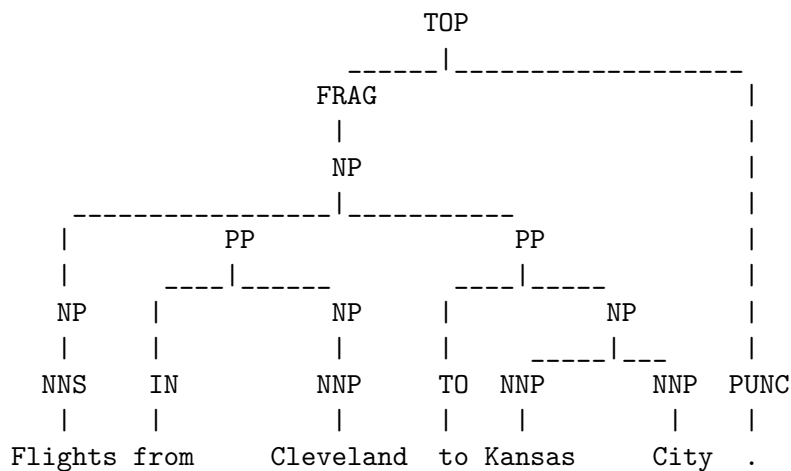
```

```
[33]: tree, logprob = cky_parse_probabilistic(grammar, tokenize('sam likes ham'))
      tree.pretty_print()
      print(f"logprob: {logprob:4.3g} | probability: {2**logprob:4.3g}")
```



logprob: -2.25 | probability: 0.21

```
[34]: # We don't use our tokenizer because the gold trees do not lowercase tokens
sent = "Flights from Cleveland to Kansas City ".split()
tree, logprob = cky_parse_probabilistic(pgrammar, sent)
tree.un_chomsky_normal_form()
tree.pretty_print()
print(f"logprob: {logprob:4.3g} | probability: {2**logprob:4.3g}")
```



logprob: -27 | probability: 7.42e-09

5.3 Evaluation of the grammar

There are a number of ways to evaluate parsing algorithms. In this project segment, you will use the “industry-standard” `evalb` implementation for computing constituent precision, recall, and F1 scores. We downloaded `evalb` during setup.

We read in the test data...

```
[35]: with open('data/test.trees') as file:
      test_trees = [Tree.fromstring(line.strip()) for line in file.readlines()]
```

```
test_sents = [tree.leaves() for tree in test_trees]
```

...and parse the test sentences using your probabilistic CKY implementation, writing the output trees to a file.

```
[36]: trees_out = []
      for sent in tqdm(test_sents):
          tree, prob = cky_parse_probabilistic(pgrammar, sent)
          if tree is not None:
              tree.un_chomsky_normal_form()
              trees_out.append(tree.pformat(margin=999999999))
          else:
              trees_out.append('()')

      with open('data/outp.trees', 'w') as file:
          for line in trees_out:
              file.write(line + '\n')
```

```
12%|          | 7/58 [00:00<00:02, 17.29it/s]100%|      | 58/58
[00:03<00:00, 16.77it/s]
```

Now we can compare the predicted trees to the ground truth trees, using `evalb`. You should expect to achieve F1 of about 0.83.

```
[37]: shell("python scripts/evalb.py data/outp.trees data/test.trees")
```

```
data/outp.trees 345 brackets
data/test.trees 471 brackets
matching        339 brackets
precision       0.9826086956521739
recall  0.7197452229299363
F1       0.8308823529411764
```

5.4 Debrief

Question: We're interested in any thoughts you have about this project segment so that we can improve it for later years, and to inform later segments for this year. Please list any issues that arose or comments you have to improve the project segment. Useful things to comment on might include the following:

- Was the project segment clear or unclear? Which portions?
- Were the readings appropriate background for the project segment?
- Are there additions or changes you think would make the project segment better? **BEGIN QUESTION** **name:** `open_response_debrief` **manual:** `true`

but you should comment on whatever aspects you found especially positive or negative.

Please stop letting students implement functions from libraries that they can simply copy-paste from their implementations. (for example, `induce_pcfg`)

6 Instructions for submission of the project segment

This project segment should be submitted to Gradescope at <https://rebrand.ly/project3-submit-code> and <https://rebrand.ly/project3-submit-pdf>, which will be made available some time before the due date.

Project segment notebooks are manually graded, not autograded using otter as labs are. (Otter is used within project segment notebooks to synchronize distribution and solution code however.)

We will not run your notebook before grading it. Instead, we ask that you submit the already freshly run notebook. The best method is to “restart kernel and run all cells”, allowing time for all cells to be run to completion. You should submit your code to Gradescope at the code submission assignment at <https://rebrand.ly/project3-submit-code>. Make sure that you are also submitting your `data/grammar` file as part of your solution code as well.

We also request that you **submit a PDF of the freshly run notebook**. The simplest method is to use “Export notebook to PDF”, which will render the notebook to PDF via LaTeX. If that doesn’t work, the method that seems to be most reliable is to export the notebook as HTML (if you are using Jupyter Notebook, you can do so using `File -> Print Preview`), open the HTML in a browser, and print it to a file. Then make sure to add the file to your git commit. Please name the file the same name as this notebook, but with a `.pdf` extension. (Conveniently, the methods just described will use that name by default.) You can then perform a git commit and push and submit the commit to Gradescope at <https://rebrand.ly/project3-submit-pdf>.

7 End of project segment 3