

ANLP Assignment 2

Ida SZUBERT
s0907677

Yue YU
s1563228

November 11, 2015

1 Tokenization

Tokenizing a sentence by identifying strings which match

| | |
|--------------------------|---|
| <code>[a-zA-Z]+</code> | tokens consisting of any number of letters |
| <code>'[a-z]+</code> | tokens consisting of ' proceeding any number of letters |
| <code>[.,?;:() -]</code> | tokens consisting of single punctuation marks |

would be sufficient to account for the most common English uses of punctuation symbols, including recognizing as separate tokens possessive 's and contractions such as 'd, 't, 're. However, there is a number of problematic cases:

- lone apostrophes as possessive suffix on nominals ending in s (such as *flowers' pollen* meaning pollen of many flowers);
- ellipsis of the centuries part of a year (such as *the '70s, in the '98*);
- informal contractions (such as *about* → *'bout*, *unless* → *'less*, *because* → *'cause*, *and* → *'n'*);
- surnames and lexical items with apostrophes (such as *'o clock*, *O'Brian*, *M'Gregor*);

The tokenize function changes the lone-apostrophes possessive into 's so that a word can be tokenized into its main lexical part and the possessive suffix. It also recognizes contracted dates. The other two problems, being largely lexical, remain unsolved.

```
def tokenise(tokenstring):
    '''Split a string into a list of tokens, treating punctuation as
    separate tokens, and splitting contractions into their parts.
    So for example "I'm leaving." --> ["I","'m","leaving","."]'''
    # normalize possessive ' to possessive 's
    normalized_1 = re.sub(r"s'\b", "s's ", tokenstring)
    # if ' is followed by more than 2 digits, it's not a year
    # split into ' and digit string
    normalized_2 = re.sub(r"'(\d{3})", r"' \1", normalized_1)
```

```

# find all alphanumeric strings, 'digit-digit strings, 's/'d/'m etc. strings
# and single punctuation marks
return re.findall(r"[a-zA-Z\d]+|'\d\d\b|'[a-z]*|[,.\?;:()-]", normalized_2)

```

2 Parse trees

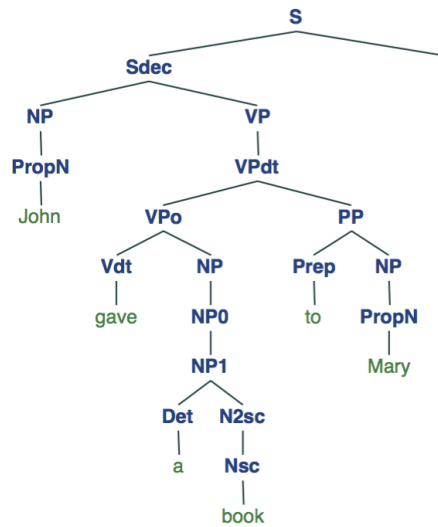


Figure 1: *John gave a book to Mary.*

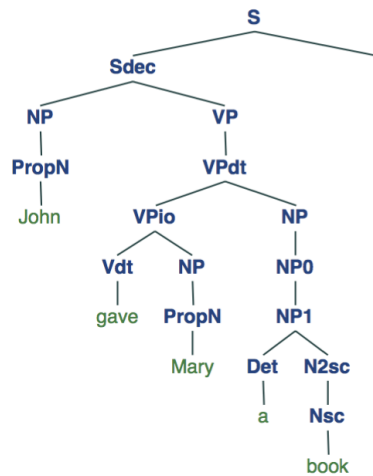


Figure 2: *John gave Mary a book.*

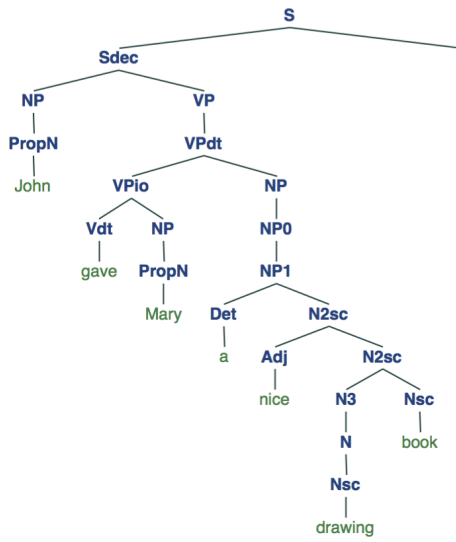


Figure 3: *John gave Mary a nice drawing book.*

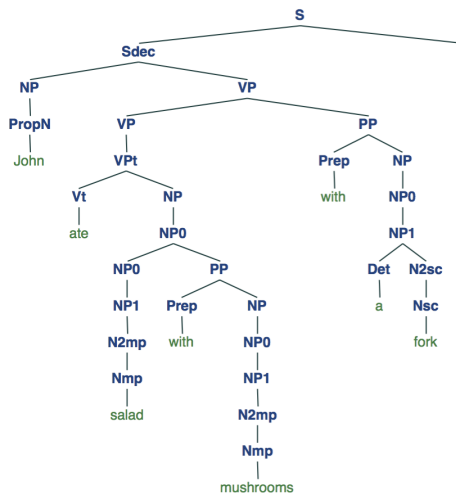


Figure 4: *John ate salad with mushrooms with a fork.*

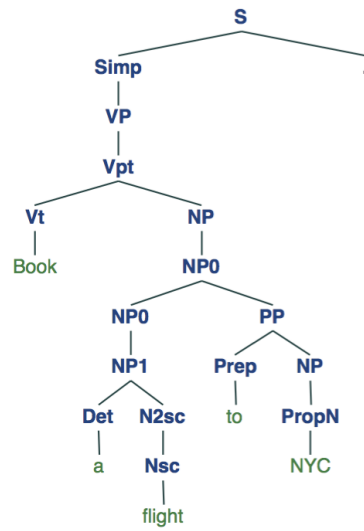


Figure 5: *Book a flight to NYC.*

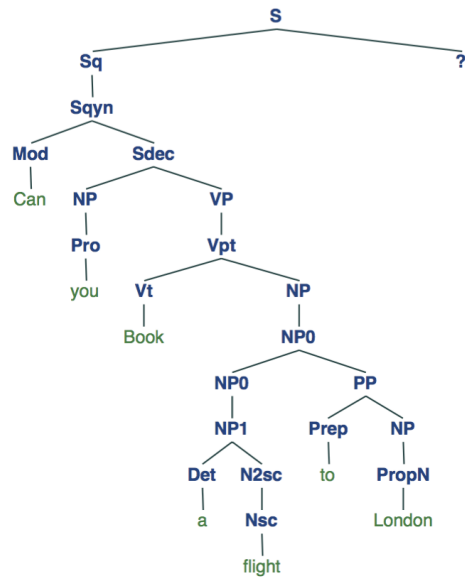


Figure 6: *Can you book a flight to London?*

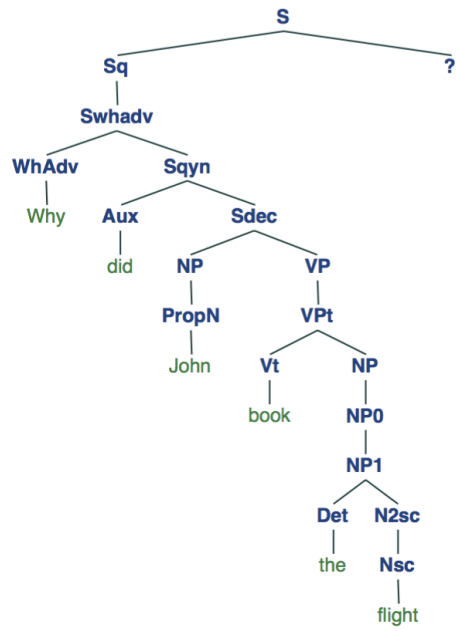


Figure 7: *Why did John book a the flight?*

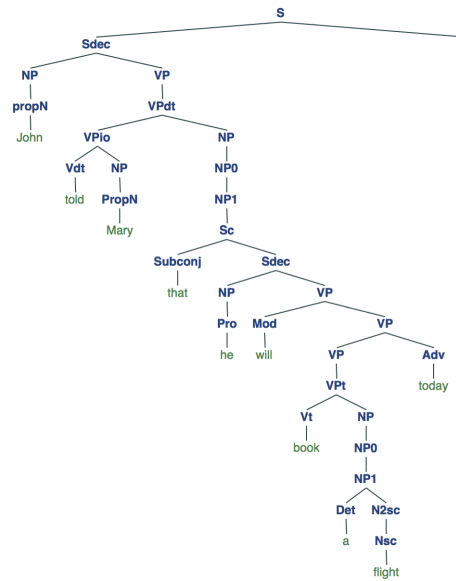


Figure 8: *John told Mary that he will book a flight today.*

3 Remarks on the grammar

3.1 Design of the grammar rules

a. Redundancy of rules in the verbal domain

There is a redundancy in the treatment of ditransitive verb phrases. VPo and $VPio$ are two nodes which expand in the exact same way:

$$VPo \rightarrow Vdt NP \text{ and } VPio \rightarrow Vdt NP.$$

and therefore there is no difference between them within the grammar. From a linguistic point of view this separation might have a motivation. The ditransitive verb and its direct object are represented by Po , while $VPio$ is a unit consisting of a ditransitive verb and its indirect object. As valid linguistic distinction as it might be, it does not need to be included in simple grammar like the one discussed here. In particular, no use is made of the distinction between that would limit overgeneration. If that was the case, then having both non-terminals would be a valid design choice. For instance, one might consider putting some restrictions on what kinds of NP can be direct or indirect objects. The needless distinction between VPo and $VPio$ propagates up to the topmost VP expansion rules. As a result, there are more VP-related rules than necessary. The grammar contains

$$\begin{aligned} VP &\rightarrow VPi \mid VPt \mid VPdt \mid Mod VP \mid VP Adv \mid VP PP \\ VPdt &\rightarrow VPo PP \\ VPdt &\rightarrow VPio NP \\ VPo &\rightarrow Vdt NP \\ VPio &\rightarrow Vdt NP \end{aligned}$$

but it could have contained as little as three rules instead:

$$\begin{aligned} VP &\rightarrow VPi \mid VPt \mid VPdt \mid Mod VP \mid VP Adv \mid VP PP \\ VPdt &\rightarrow VPd NP \mid VPd PP \\ VPd &\rightarrow Vdt NP \end{aligned}$$

b. Redundancy of rules in the nominal domain

Redundancy can be also observed among the nominal rules. The NP_0 non-terminal node seems redundant. Its purpose is to allow for recursive generation of prepositional phrases following a noun phrase. This can be achieved by accounting for this kind of recursion within one of the existing NP rules. We could propose two sets of rules to replace the following:

$$\begin{aligned}
NP &\rightarrow \text{PropN} \mid \text{Pro} \mid NP_0 \\
NP_0 &\rightarrow NP_1 \mid NP_0 PP \\
NP_1 &\rightarrow \text{Det} \text{ N2sc} \mid \text{N2mp} \mid \text{Sc}
\end{aligned}$$

The first replacement option moves recursive *PP* addition into an *NP* expansion rule:

$$\begin{aligned}
NP &\rightarrow \text{PropN} \mid \text{Pro} \mid NP_1 \mid NP PP \\
NP_1 &\rightarrow \text{Det} \text{ N2sc} \mid \text{N2mp} \mid \text{Sc}
\end{aligned}$$

Such rearrangement has the advantage of allowing for not only common nouns, but also proper names and pronouns to be modified by *PPs*. This could be desirable in a more comprehensive grammar, but additional restrictions would be required, e.g. to account for pronouns generally not taking *PP* complements. The second option moves the recursion into an *NP₁* expansion:

$$\begin{aligned}
NP &\rightarrow \text{PropN} \mid \text{Pro} \mid NP_1 \\
NP_1 &\rightarrow \text{Det} \text{ N2sc} \mid \text{N2mp} \mid \text{Sc} \mid NP_1 PP
\end{aligned}$$

It is more conservative in that with these rules in replacing the originals, the grammar would generate the same strings, i.e. the grammars with and without replacement would be weakly equivalent [1]. The *NP₀* node expands to an *NP₁* followed by an arbitrary number of *PPs*. The same effect will be achieved by recursively expanding *NP₁* into an *NP₁* and *PP*, with the benefit of reducing the number of non-terminals.

One might argue that linguistic justification of the existence of the *NP₀* can be found in its correspondence with N' node in X-bar syntax. Both are projections above N and below full NP, and allow for adding in principle an unlimited number of modifiers to the head [2]. However, the unlimited recursive expansion of the *NP₀* node implies that a noun can have many PP adjuncts, which does not seem justified. Nominal PP adjuncts do exist, for instance in the NP *the teapot on the table*, but whether they are stackable is not obvious [3]. We can say *the teapot on the table in the kitchen*, but the second PP might be modifying *table* rather than *teapot*. This is clear in the phrase *the teapot on the table with a crack*, where it's the table that has a crack, not the teapot. The similarity between *NP₀* and N' thus seems superficial and can be rejected as motivation for including the former in the grammar.

c. **Dobtfull usefulness of the *Sq* node**

The *Sq* node serves no purpose other than to reflect the conceptual relation between *Sqyn* and *Swhadv*, namely the fact both are non-terminals expanding into questions. This, however, represents only a superficial,

if not mistaken, grammatical insight. In fact wh-questions and auxiliary-inversion questions are represented by different structures (and are created by different kinds of movement processes [4]). If the motivation for including the *Sq* node is linguistic, then the reasons seem insufficient. If, on the other hand, the node was included in the grammar to simplify the expansion of *S*, then by analogy we should also include

$$Sd \rightarrow Sdec \mid Simp$$

and change the expansion of *S* to

$$S \rightarrow Sd \text{ '}' \mid Sq \text{ '}'?$$

d. **Lower case and capitalized versions of terminals.**

It does not seem necessary to include lower case and capitalized versions of the same terminals. In the sample sentences we have seen *Book*, *Can*, and *Why* being capitalized to account for their sentence-initial positions. Instead of duplicating lexical items, we could simply capitalize the first word of a generated sentence and ignore sentence-initial capitalization during parsing. There would be no ill effects of such a change, since every generated sentence is, by definition, grammatical. Every word which happens to be initial in a sentence produced by the grammar is allowed to be there. The additional confirmation of this fact in form of the word's capitalised version being present in the grammar is unnecessary.

3.2 Overgeneration

a. **Lack of subject-verb agreement**

There's distinction between single/count and plural/mass nouns, but no corresponding distinction in the verbal domain. The problem is somewhat masked by the fact that most of the included verbs are in past tense, where in English there are no subject-verb agreement phenomena. Nevertheless, we can see that the grammar would accept **the fork book a flight.*, a sentence which ignores verb conjugation requirements.

b. **Treating VP as full sentences**

The rules:

$$S \rightarrow Simp \text{ '}'$$

$$Simp \rightarrow VP$$

imply that any *VP* can be a grammatical sentence on its own. This, however, is true only of *VP* in imperative mood. The grammar, however, allows for treating any *VP*, e.g. **ate today*, as a full sentence. Solving this problem is not possible without substantially expanding the grammar. We

would have to distinguish between infinitival and inflected forms of verbs, and extend the set of possible expansions of *VP* accordingly. We could then state that *Simp* subcategories for *VPs* whose head is an infinitival form of a verb. Therefore, we would need the following rules:

$$\begin{aligned}
Vt_{inf} &\rightarrow \text{'Book'} \mid \text{'Tell'} \mid \dots \\
VPi_{inf} &\rightarrow \text{'Eat'} \mid \dots \\
VPt_{inf} &\rightarrow Vt_{inf} \text{ NP} \mid VPt_{inf} \text{ Adv} \mid VPt_{inf} \text{ PP} \\
VP &\rightarrow VPt_{inf} \\
\text{Simp} &\rightarrow VPt_{inf} \mid VPi_{inf}
\end{aligned}$$

Even then, additional restrictions would be needed to account for characteristics of the imperative mood, e.g. the fact that the second person pronoun needs to be in its reflexive form when used as a direct object of the matrix verb [5]. Otherwise the grammar would generate sentences such as **Book you a flight to London.* instead of grammatical **Book yourself a flight to London.*

c. **Case agreement**

The grammar does not account for case subcategorisation requirements of verbs. The only corner of the modern English syntax which still exhibits case phenomena is the pronoun system [6], however the grammar does not include a *ProNom* and a *ProAcc* pre-terminals. Having no distinction between nominative and accusative forms leads to admitting non-grammatical sentences, e.g. **Mary gave **he** a nice salad.* instead of *Mary gave **him** a nice salad.*

d. **Verb and noun subcategorization for pronouns**

Verbs and nouns tend to co-occur with particular pronouns more often than with others. However, the grammar does not account for this preference. For instance, one can *eat with a fork*, *eat with Mary*, or have a *flight to London*, but the grammar also allows for having a **salad to fork* and booking a **flight with NYC*.

e. **Multiple modals**

The grammar allows for an unlimited number of modals to precede a *VP*. This can lead to generating sentences such as **Can will can can John book a flight?*. However, English limits the acceptable number (up to 4, depending on dialect), orderings (evidential, epistemic, deontic), and identity (*might* or *may* being the most commonly first in the sequence (Di Paolo 1989)) of verbs in multiple modal constructions. In order to capture these restrictions we would need a whole series of serially connected expansion rules, each introducing one modal of a certain type.

4 Comments for CKY.buildIndices

Answer to be found in cky.py

5 Comments for CKY.unary_fill & Cell.unary_update

Answer to be found in cky.py

6 Comments for CKY.parse & CKY.maybe_build

Answer to be found in cky.py

7 Ambiguity Analysis

7.1 "John gave a book to Mary."

Distinct Parse(s): 3

Ambiguity Source: 3 different **VPs** found in (1,6) can pair up with the **NP** in (0,1) to form **Sdecl**. The **VPs** came from the following rules:

$$\begin{aligned} VP &\rightarrow VP(1,4) PP(4,6) \\ VP &\rightarrow VPdt \rightarrow VPo(1,4) PP(4,6) \\ VP &\rightarrow VPt \rightarrow Vt(1,2) NP(2,6) \end{aligned}$$

Interpretation of ambiguity:

[VP [VP John gave a book] [PP to Mary]]

In this parse, *John gave a book* is a full VP, and *to Mary* is an additional modifier. This analysis would be correct for a sentence such as *John read a book on the couch*, where *the couch* is not an argument of the verb. In case of *give* both the NP and the PP are arguments.

[VP [VPdt [VPo John gave a book] [PP to Mary]]]

This is the most appropriate parse, in which *give* has a direct object, *a book*, and an indirect one, *Mary*. The only VP in this analysis is the one which correctly encompasses all of the participants of the event.

[VP [VPt [Vt John gave] [NP a book to Mary]]]

This analysis incorrectly treats *to Mary* as forming as a modifier of *a book*. This structure would correctly describe a sentence such as *John ate pasta with tuna*, where the PP really modifies the noun, and the verb really is transitive.

7.2 "John gave Mary a book."

Distinct Parse(s): 1

7.3 "John gave Mary a nice drawing book."

Distinct Parse(s): 2

Ambiguity Source: 2 different **N2sc** found in (4,7) can pair up with the **Det** in (3,4) to form **NP1**. The **N2sc** came from the following rules:

$$\begin{aligned} N2sc &\rightarrow Adj(4,5) N2sc(5,7) \\ N2sc(5,7) &\rightarrow N3(5,6) Nsc(6,7) \\ N2sc(5,7) &\rightarrow Adj(5,6) N2sc(6,7) \end{aligned}$$

Interpretation of ambiguity:

[N2sc [Adj drawing] [Nsc book]]

The grammar lists *drawing* as an adjective, so the unit *drawing book* can be analysed as a noun modified by an adjective, in the way *fragile* is a modifier in *fragile cup*.

[N2sc [N3 [N drawing]] [Nsc book]]

In this parse both *drawing* and *book* are treated as nouns, and N2sc represents a noun-noun compound, in the way *tea cup* is a compound of two nouns. This seems to be a correct interpretation, given how popular such structures are in English, and that there is no evidence of *drawing* being an adjective (e.g. there are no comparative forms **a drawinger book*).

7.4 "John ate salad with mushrooms with a fork."

Distinct Parse(s): 5

Ambiguity Source: 5 different **VP** found in (1,8) can pair up with the **NP** in (0,1) to form **Sdecl**. The **VP** came from the following rules:

$$\begin{aligned} VP &\rightarrow VP(1,3) PP(3,8) \\ VP &\rightarrow VP(1,5) PP(5,8) \\ VP(1,5) &\rightarrow VP(1,3) PP(3,5) \\ VP(1,5) &\rightarrow VPt \rightarrow Vt(1,2) NP(2,5) \\ VP &\rightarrow VPt \rightarrow Vt(1,2) NP(2,8) \rightarrow NP0 \\ NP0 &\rightarrow NP0(2,3) PP(3,8) \\ NP0 &\rightarrow NP0(2,5) PP(5,8) \end{aligned}$$

Interpretation of ambiguity:

[VP [VP ate salad] [PP with mushrooms with a fork]]

On this interpretation the mushrooms, who have the fork, were keeping John company while he was eating the salad. It is a nonsensical parse in this case, but would be a correct one for a sentence such as *John planned a trip with a friend with disability*.

[VP [VP [VP ate salad] [PP with mushrooms]] [PP with a fork]]

In this parse *with mushrooms* is a modifier of a salad-eating event. This structure would be a good representation for a sentence such as *John ate salad with*

pleasure on a rooftop terrace where there is an event, namely salad eating, and two PP modifiers.

[VP [VP [VPt [Vt ate] [NP salad with mushrooms]]] [PP with a fork]]

This is the correct analysis. It captures the fact that *salad with mushrooms* is a constituent and that it should be treated as an object of the verb. The PP is correctly represented as a modifier of a full VP, since it is not an argument of a verb.

[VP [VPt [Vt ate] [NP [NP0 [NP0 salad] [PP with mushrooms with a fork]]]]]

In this analysis the PP *with a fork* modifies its closest noun, hence the salad is full of mushrooms, which have a fork. It would be an appropriate parse for a sentence such as *John ate a salad with truffles from Piedmont*.

[VP [VPt [Vt ate] [NP [NP0 [NP0 salad with mushrooms] [PP with a fork]]]]]

Finally, the two PPs following *salad* can be treated as two separate modifiers of the noun. On this reading the salad is with mushrooms, and it is also with a fork. Again, the interpretation is nonsensical. It could be correctly applied to the sentence *Jane admired the earrings with sapphires from her grandmother*, on the condition that the earrings were from the grandmother, not the sapphires.

7.5 "Book a flight to NYC."

Distinct Parse(s): 2

Ambiguity Source: 2 different **Simp** found in (0,5) can pair up with the '.' in (5,6) to form **S**. The **Simp** came from the following rules:

$$\begin{aligned} \text{Simp} &\rightarrow \text{VP} \\ \text{VP} &\rightarrow \text{VP}(0,3) \text{PP}(3,5) \\ \text{VP} &\rightarrow \text{VPt} \rightarrow \text{Vt}(0,1) \text{NP}(1,5) \end{aligned}$$

Interpretation of ambiguity:

[VP [VP Book [NP a flight]] [PP to NYC]]

In this parse the booking of a flight happens to NYC, which does not represent the meaning of the sentence correctly. This structure would be appropriate for a sentence such as *Book a flight with my credit card*, where the final PP is indeed a modifier of the booking event.

[VP [VPt Book [NP a flight [PP to NYC]]] This is the correct analysis, on which someone is to book a flight whose destination is NYC.

7.6 "Can you book a flight to London?"

Distinct Parse(s): 2

Ambiguity Source: 2 different **Sq** found in (0,7) can pair up with the '??' in (7,8) to form **S**. The **Sq** came from the following rules:

$$\begin{aligned} \text{Sq} &\rightarrow \text{Sqyn} \rightarrow \text{Mod}(0,1) \text{Sdecl}(1,7) \\ \text{Sdecl}(1,7) &\rightarrow \text{NP}(1,2) \text{VP}(2,7) \end{aligned}$$

$$\begin{aligned}
&VP(2,7) \rightarrow VP(2,5) PP(5,7) \\
&VP(2,7) \rightarrow VPt \rightarrow Vt(2,3) NP(3,7)
\end{aligned}$$

Interpretation of ambiguity:

[VP [VP book a flight] [PP to London]]

[VP [VPt [Vt book] [NP a flight to London]]]

This ambiguity is discussed in section 7.5

7.7 "Why did John book the flight?"

Distinct Parse(s): 1

7.8 "John told Mary that he will book a flight today."

Distinct Parse(s): 3

Ambiguity Source: 3 different **Sdecl** found in (0,10) can pair up with the '.' in (10,11) to form **S**. The **Sdecl** came from the following rules:

$$\begin{aligned}
&Sdecl \rightarrow NP(0,1) VP(1,10) \\
&VP(1,10) \rightarrow VP(1,9) Adv(9,10) \\
&VP(1,10) \rightarrow VPdt \rightarrow VPio(1,3) NP(3,10) \\
&NP(3,10) \rightarrow NP0 \rightarrow NP1 \rightarrow Sc \rightarrow Subconj(3,4) Sdecl(4,10) \\
&Sdecl(4,10) \rightarrow NP(4,5) VP(5,10) \\
&VP(5,10) \rightarrow Mod(5,6) VP(6,10) \\
&VP(5,10) \rightarrow VP(5,9) Adv(9,10)
\end{aligned}$$

Interpretation of ambiguity:

[VP [VP told Mary that he will book a flight] [Adv today]]

On this analysis the adverb modifies the telling event, not the booking event, which is a sensible interpretation of the sentence. However, it is more natural in English to move heavy constituents to the left, and so we might have expected *John told Mary today that he will book a flight*.

[VP [Mod will] [VP [VP book a flight] [Adv today]]]

The difference between the two other interpretations is subtle. In both cases the adverb modifies booking rather than telling. In the analysis above John books something today. This treats *today* as a VP adverb.

[VP [VP [Mod will] [VP book a flight] [Adv today]]]

In the alternative parse, John will do something today. This treats *today* as a sentential adverb.

The difference between two parses lies in whether *today* adjoins to a VP or an AuxP/IP (auxiliary aka inflectional phrase). Distribution of adverbs is complex and in principle an adverb might modify either of those phrases. In the case at hand the VP-adjunction analysis seems to be the correct one, since sentential adverbs generally do not appear clause-finally [7].

8 Generating a parse tree

8.1 CKY parsing

Up to now, the program is a CKY recognizer that can only determine whether a string belongs to the language generated by the given grammar. In order to extend it into a CKY parser that can construct a parse tree, back-pointers need to be attached to each label in each cell. The backpointers record the one or more ways in which a constituent spanning (i,j) can be made from constituents spanning (i,k) and (k,j) , for any k such that $i < k < j$. Therefore, having found an S node in the top-right-hand corner of the matrix, we can trace back from this node through every intermediate constituent down to the terminals which form the string. The final result is then a shared-forest of possible parse trees, where common tree parts are factored between the various parses [8].

8.2 Design and implementation

First, a Label class was defined to use in place of Strings for describing the constituent labels in each cell. Taking advantage of the flexible data structure in Python, this class can be initialized with at most 4 parameters, as shown below.

List of parameters:

- **symbol** the represented label (required)
- **start** row number of the cell (optional)
- **end** column number of the cell (optional)
- **rhs** list of right-hand-sides of the production rules which generate the target symbol; mid value is also included in the list, which refers to the source cells from which the constituents of right-hand-sides came (optional)

The Label class can take different number of parameters depending on where the label came from, such as the left-hand-sides of rules and words from the input string. For instance, if a symbol was added according to a unary rule, the row number and column number will be the same as those of the parent label. For the labels of words, which will be the leaf nodes of a parse tree, the value of right-hand-side will be omitted (not available).

| Source Type Parameter | binary rule | unary rule | word |
|--------------------------|-------------|------------|------|
| symbol | ✓ | ✓ | ✓ |
| start | ✓ | | |
| end | ✓ | | |
| rhs | ✓ | ✓ | |

Table 1: Parameter Structure for Different Labels

Moreover, a number of the existing CKY and Cell methods were edited to construct/exploit this richer label structure. Three utility functions for listing, searching and updating label list are also provided.

After that, the backtrack algorithm was implemented in two functions:

- **update_trees** to update all the possible subtrees of the give label in a cell
(In this implementation, the function will be ended after the first subtree was returned)
- **update_children** to construct NLTK tree nodes according to given rules

update_trees and *update_children* will call each other recursively to generate the whole parse tree.

8.3 Summary

To generate the first parse tree, we only need to pass the Start Symbols (S) in the top-right-hand corner of the matrix to *update_trees* method and it will return the result after the first parse tree was found. However, in order to return all the possible parse trees, the Start Symbols (S) in the top-right-hand corner of the matrix should be treated differently from other labels and *update_trees* method should return a list of possible subtrees to construct all the ambiguous result. The detail of generating all parse trees will not be discussed in this report.

References

- [1] D. Jurafsky and J. Martin, *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Upper Saddle River, N.J: Prentice Hall, 2008.
- [2] G. Kornai, Andras; Pullum, “The x-bar theory of phrase structure,” *Language*, pp. 24–50, 1990.
- [3] N. Hedberg and R. DeArmond, “On nominal complements and adjuncts,” *Proceedings of the 1999 Annual Conference of the Canadian Linguistics Association*, 1999.
- [4] A. Carnie, *Syntax. A generative introduction*. Oxford, England: Blackwell, 2013.
- [5] F. Palmer, *Mood and modality*. Cambridge: Cambridge University Press, 2001.
- [6] R. Huddleston and G. Pullum, *The Cambridge Grammar of the English Language*. Cambridge: Cambridge University Press, 2002.
- [7] E. Potsdam, “A syntax for adverbs,” *WECOL 98, Proceedings of the Western Conference on Linguistics*, pp. 397–411, 1998.
- [8] B. Lang, “Recognition can be harder than parsing,” *Computational Intelligence*, vol. 10, no. 4, pp. 486–494, 1994.