# Handout 15: Parsing

## NLTK chart parser

1. A toy grammar

    a. Create a file called `g1.cfg` containing:

    ```
    % start S

    S -> NP VP
    VP -> V NP | V NP PP
    PP -> P NP
    NP -> Det N | Det N PP | Name

    V -> "saw" | "ate" | "walked"
    Name -> "John" | "Mary" | "Bob"
    Det -> "a" | "an" | "the" | "my"
    N -> "man" | "dog" | "cat" | "telescope" | "park"
    P -> "in" | "on" | "by" | "with"
    ```

    b. In Python:

    ```
    >>> s = open('g1.cfg').read()
    >>> from nltk import CFG
    >>> g = CFG.fromstring(s)
    >>> type(g)
    <class 'nltk.grammar.CFG'>
    >>> print(g)
    Grammar with 26 productions (start state = S)
        S -> NP VP
        VP -> V NP
        ...
    ```

2. Create and use a parser

    ```
    >>> from nltk import ChartParser
    >>> p = ChartParser(g)
    >>> s = 'the dog saw the cat'.split()
    >>> s
    ['the', 'dog', 'saw', 'the', 'cat']
    >>> for t in p.parse(s):
    ...     print(t)
    ...
    (S (NP (Det the) (N dog)) (VP (V saw) (NP (Det the) (N cat))))
    ```

**3.** That's a lot of steps to remember. Assignment: package them up as a
function called `make_parser`:

```
1    >>> p = make_parser('g1.cfg')
2    >>> print(next(p.parse(s)))
3    (S (NP (Det the) (N dog)) (VP (V saw) (NP (Det the) (N cat))))
```

**4.** Note: if you have a parser p, you can get the grammar:

```
1    >>> p.grammar()
2    <Grammar with 26 productions>
```

**5.** What if the sentence fails to parse?

    **a.** Not very useful:

```
1    >>> s = 'the dog with the telescope walked'.split()
2    >>> p.parse(s).next()
3    Traceback (most recent call last):
4      File "<stdin>", line 1, in <module>
5    StopIteration
```

    **b.** Getting a trace:

```
1    >>> p.chart_parse(s, trace=2)
2    |. the   . dog   . with . the   .telesc.walked.|
3    ...
4    Single Edge Fundamental Rule:
5    |[---------------------------------]      .| [0:5] NP -> Det N PP *
6    ...
7    Bottom Up Predict Combine Rule:
8    |.       .       .       .       .       [------>| [5:6] VP -> V * NP
9    |.       .       .       .       .       [------>| [5:6] VP -> V * NP PP
10   <nltk.parse.chart.Chart object at 0x101c9f850>
```

# Trees

**6.** Anatomy of an NLTK tree

    **a.** Creating a tree manually

```
1    >>> from nltk import Tree
2    >>> t = Tree.fromstring('(S (NP Fido) (V barked))')
3    >>> print(t)
4    (S (NP Fido) (V barked))
```

    **b.** The structure

```
1    >>> t
2    Tree('S', [Tree('NP', ['Fido']), Tree('V', ['barked'])])
```

7. Tree = node

    a. Label

```
1    >>> t.label()
2    'S'
```

    b. The children of the node: tree behaves like a list

```
1    >>> len(t)
2    2
3    >>> print(t[0])
4    (NP Fido)
5    >>> print(t[1])
6    (V barked)
```

8. Leaf nodes

    a. The leaves are strings, not trees

```
1    >>> np = t[0]
2    >>> print(np)
3    (NP Fido)
4    >>> len(np)
5    1
6    >>> n = np[0]
7    >>> n
8    'Fido'
```

    b. Testing for a leaf

```
1    >>> isinstance(np, str)
2    False
3    >>> isinstance(n, str)
4    True
```

9. Creating a tree

```
1    >>> subj = Tree('NP', ['Garfield'])
2    >>> obj = Tree('NP', ['lasagna'])
3    >>> v = Tree('V', ['likes'])
4    >>> vp = Tree('VP', [v, obj])
5    >>> s = Tree('S', [subj, vp])
6    >>> print(s)
7    (S (NP Garfield) (VP (V likes) (NP lasagna)))
```

10. Functions on trees

    a. Comparison of trees

```
1    >>> t == Tree.fromstring('(S (NP Spot) (V ran))')
2    False
3    >>> t == Tree.fromstring('(S (NP Fido) (V barked))')
4    True
```

**b.** Subtrees

```
1    >>> for st in t.subtrees(): print(st)
2    ...
3    (S (NP Fido) (V barked))
4    (NP Fido)
5    (V barked)
```
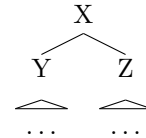
**c.** Leaves

```
1    >>> t.leaves()
2    ['Fido', 'barked']
3    >>> t.pos()
4    [('Fido', 'NP'), ('barked', 'V')]
```

**11.** Defining recursive functions on tree structure

    **a.** Write a function that counts the nodes in a tree. For example:

```
1    >>> count_nodes(t)
2    3
3    >>> count_nodes(np)
4    1
```

    **b.** Answer: how many nodes in a tree? One for $X$, plus the number in the $Y$ subtree, plus the number in the $Z$ subtree.

```
1    n = 1
2    for child in tree:
3        n += count_nodes(child)
```

    **c.** "Bottoming out": when you get to a terminal node (string), the answer is 0 (not a node!).

```
1    def count_nodes (tree):
2        if isinstance(tree, str):
3            return 0
4        else:
5            n = 1
6            for child in tree:
7                n += count_nodes(child)
8            return n
```

**12.** Another example: make a mirror image of a tree in which all expansions are reversed from the original.

13. Saving trees

    a. Load `g1.cfg`:

```
1    >>> p = make_parser('g1.cfg')
2    >>> s = 'the dog saw the cat'.split()
3    >>> trees = list(p.parse(s))
```

    b. Converting a tree to a string:

```
1    >>> str(trees[0])
2    '(S (NP (Det the) (N dog)) (VP (V saw) (NP (Det the) (N cat))))'
```

    c. We know how to write a file

```
1    >>> with open('sents1.txt', 'w') as f:
2    ...      f.write(str(trees[0]) + '\n')
3    ...
```

14. Loading trees

    a. Reload

```
1    >>> f = open('sents1.txt')
2    >>> line = next(f).rstrip('\r\n')
3    >>> t = Tree.fromstring(line)
4    >>> print(t)
5    (S (NP (Det the) (N dog)) (VP (V saw) (NP (Det the) (N cat))))
```

    b. Suppose `g1.cfg` has been modified. Are we still getting the same
    result as before?

```
1    >>> p = make_parser('g1.cfg')
2    >>> trees = list(p.parse(s))
3    >>> t == trees[0]
4    True
```

15. Penn treebank (sample)

```
1    >>> from nltk.corpus import treebank
2    >>> trees = treebank.parsed_sents()
3    >>> print(trees[0])
4    (S
5      (NP-SBJ
6        (NP (NNP Pierre) (NNP Vinken))
7          ...
```