

# Induction of probabilistic context-free grammars from tree corpora

Praktikum Parsing von natürlichen Sprachen

Tobias Denking

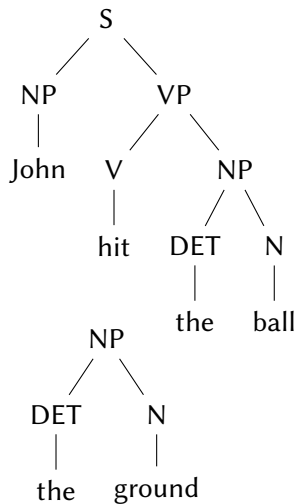
2019-04-09

# Overview

- 1 Idea (recap)
- 2 Overall pipeline (suggestion)
- 3 Trees and s-expressions
- 4 Data structures (suggestion)
- 5 General advice

## Idea (recap)

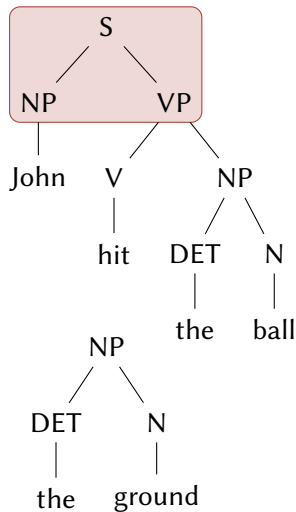
rules:



## Idea (recap)

rules:

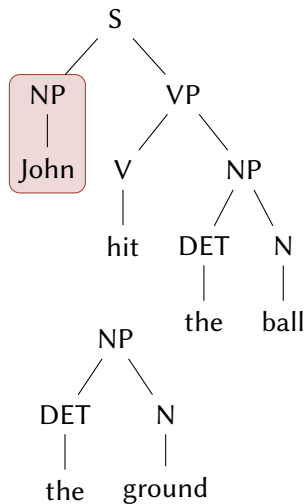
- $S \rightarrow NP VP$



## Idea (recap)

rules:

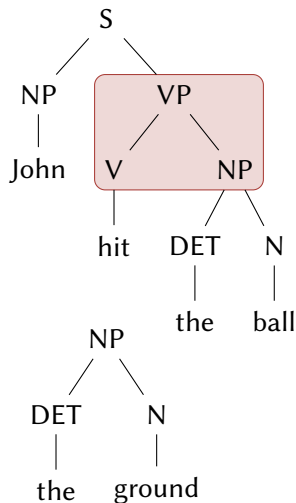
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$



## Idea (recap)

rules:

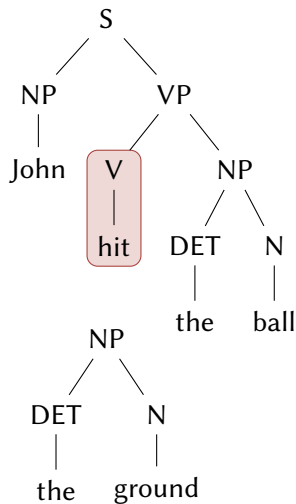
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$



## Idea (recap)

rules:

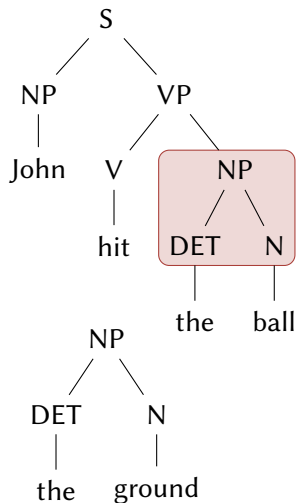
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$



## Idea (recap)

rules:

- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$
- $NP \rightarrow \text{DET } N$

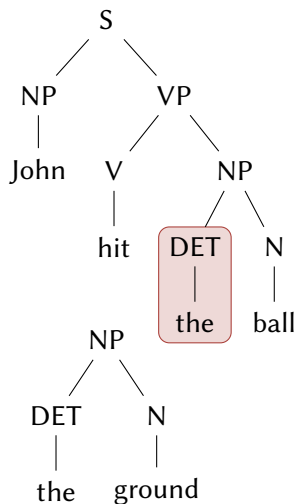




## Idea (recap)

rules:

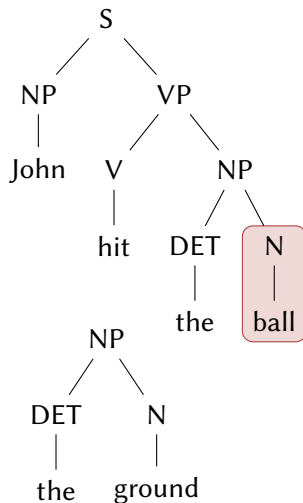
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$
- $NP \rightarrow \text{DET N}$
- $\text{DET} \rightarrow \text{the}$



## Idea (recap)

rules:

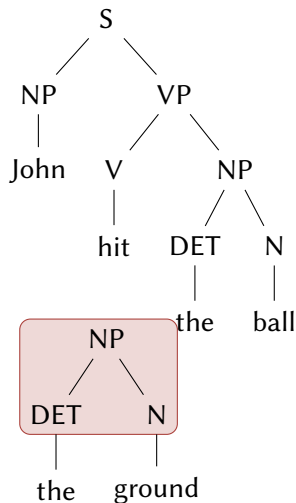
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$
- $NP \rightarrow \text{DET } N$
- $\text{DET} \rightarrow \text{the}$
- $N \rightarrow \text{ball}$



## Idea (recap)

rules:

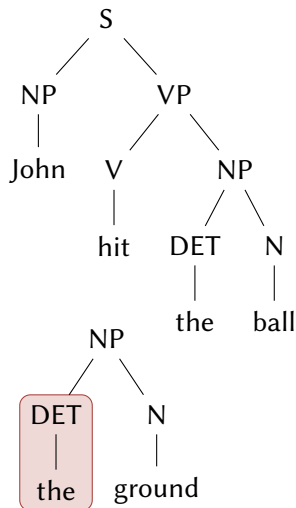
- $S \rightarrow NP VP$  # 1
- $NP \rightarrow John$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1
- $NP \rightarrow DET N$  # 2
- $DET \rightarrow the$  # 1
- $N \rightarrow ball$  # 1



# Idea (recap)

rules:

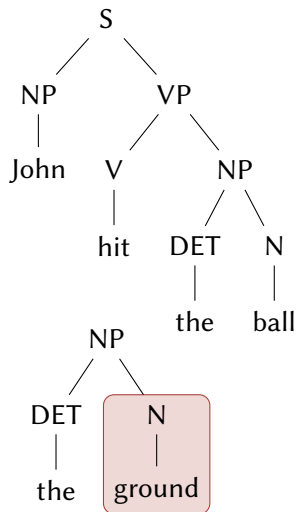
- $S \rightarrow NP VP$  # 1
- $NP \rightarrow John$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1
- $NP \rightarrow DET N$  # 2
- $DET \rightarrow the$  # 2
- $N \rightarrow ball$  # 1



## Idea (recap)

rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow \text{John}$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow \text{hit}$  # 1
- $NP \rightarrow \text{DET N}$  # 2
- $\text{DET} \rightarrow \text{the}$  # 2
- $N \rightarrow \text{ball}$  # 1
- $N \rightarrow \text{ground}$  # 1



## Idea (recap)

rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow John$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1
- $NP \rightarrow DET N$  # 2
- $DET \rightarrow the$  # 2
- $N \rightarrow ball$  # 1
- $N \rightarrow ground$  # 1

normalised rules:

- $S \rightarrow NP VP$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1

## Idea (recap)

rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow \text{John}$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow \text{hit}$  # 1
- $NP \rightarrow \text{DET N}$  # 2
- $\text{DET} \rightarrow \text{the}$  # 2
- $N \rightarrow \text{ball}$  # 1
- $N \rightarrow \text{ground}$  # 1

normalised rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow \text{John}$  #  $\frac{1}{1+2} = 1/3$
- $VP \rightarrow V NP$  # 1
- $V \rightarrow \text{hit}$  # 1
- $NP \rightarrow \text{DET N}$  #  $\frac{2}{1+2} = 2/3$

## Idea (recap)

rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow John$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1
- $NP \rightarrow DET N$  # 2
- $DET \rightarrow the$  # 2
- $N \rightarrow ball$  # 1
- $N \rightarrow ground$  # 1

normalised rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow John$  #  $\frac{1}{1+2} = 1/3$
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1
- $NP \rightarrow DET N$  #  $\frac{2}{1+2} = 2/3$
- $DET \rightarrow the$  #  $\frac{2}{2} = 1$



## Idea (recap)

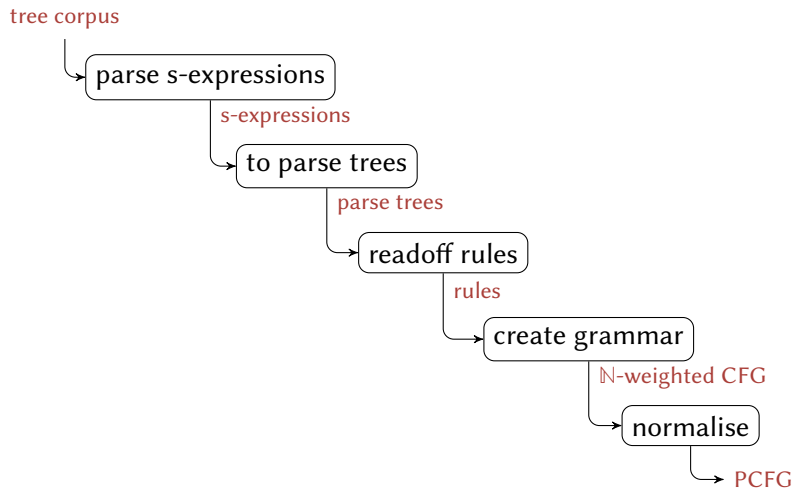
rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow John$  # 1
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1
- $NP \rightarrow DET N$  # 2
- $DET \rightarrow the$  # 2
- $N \rightarrow ball$  # 1
- $N \rightarrow ground$  # 1

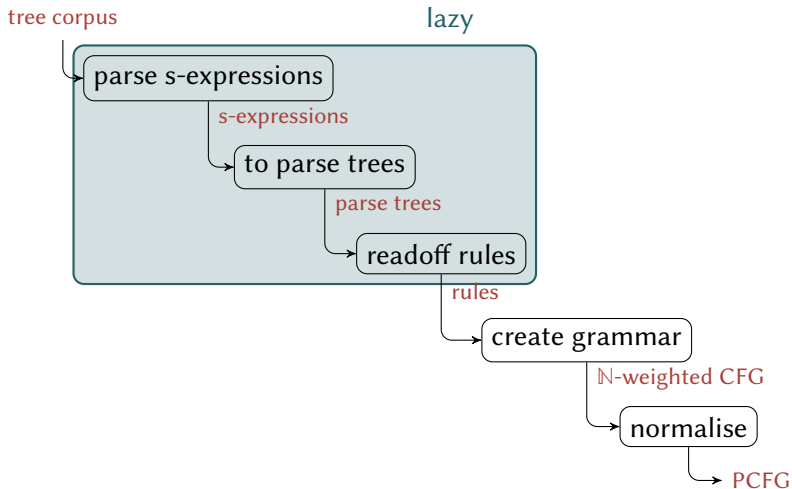
normalised rules:

- $S \rightarrow NP VP$  # 1
- $NP \rightarrow John$  #  $\frac{1}{1+2} = 1/3$
- $VP \rightarrow V NP$  # 1
- $V \rightarrow hit$  # 1
- $NP \rightarrow DET N$  #  $\frac{2}{1+2} = 2/3$
- $DET \rightarrow the$  #  $\frac{2}{2} = 1$
- $N \rightarrow ball$  #  $\frac{1}{1+1} = 1/2$
- $N \rightarrow ground$  #  $\frac{1}{1+1} = 1/2$

# Overall pipeline (suggestion)



# Overall pipeline (suggestion)



# Trees and s-expressions

- *s-expression over a set  $A$ :*

symbolic **expressions**

**atoms:** elements of  $A$

**lists:** strings  $(s_1 \sqcup s_2 \sqcup \dots \sqcup s_k)$  where the  $s_i$ 's are s-expressions

# Trees and s-expressions

- *s-expression over a set  $A$ :* symbolic **expressions**
  - atoms: elements of  $A$
  - lists: strings  $(s_1 \sqcup s_2 \sqcup \dots \sqcup s_k)$  where the  $s_i$ 's are s-expressions
- data structure in Rust (using tagged union):

```
1 pub enum SExp<A> {  
2     Atom(A),  
3     List(Vec<SExp<A>>),  
4 }
```

# Trees and s-expressions

- *s-expression over a set  $A$ :* symbolic expressions
  - atoms: elements of  $A$
  - lists: strings  $(s_1 \sqcup s_2 \sqcup \dots \sqcup s_k)$  where the  $s_i$ 's are s-expressions
- data structure in Rust (using tagged union):

```
1 pub enum SExp<A> {  
2     Atom(A),  
3     List(Vec<SExp<A>>),  
4 }
```

- trees as s-expressions:
  - no empty lists
  - $s_1$  is always in  $A$

# Data structures (suggestion)

- tree

```
1 pub struct Tree<A> {  
2   pub root: A,  
3   pub children: Vec<Tree<A>>  
4 }
```

# Data structures (suggestion)

- tree

```
1 pub struct Tree<A> {  
2     pub root: A,  
3     pub children: Vec<Tree<A>>  
4 }
```

- rule

```
1 pub enum Rule<N, T> {  
2     Lexical { lhs: N, rhs: T, },  
3     NonLexical { lhs: N, rhs: Vec<N>, },  
4 }
```



# Data structures (suggestion)

- tree

```
1 pub struct Tree<A> {  
2     pub root: A,  
3     pub children: Vec<Tree<A>>  
4 }
```

- rule

```
1 pub enum Rule<N, T> {  
2     Lexical { lhs: N, rhs: T, },  
3     NonLexical { lhs: N, rhs: Vec<N>, },  
4 }
```

- grammar

```
1 pub struct Grammar<N, T, W> where N: Eq + Hash, T: Eq + Hash {  
2     initial: N,  
3     rules: HashMap<Rule<N, T>, W>,  
4 }
```

# General advice

- use *unit tests*

# General advice

- use *unit tests*
- use a *build system*, e.g. cargo (Rust), stack or cabal (Haskell)

# General advice

- use *unit tests*
- use a *build system*, e.g. cargo (Rust), stack or cabal (Haskell)
- use a *cli library*, e.g. clap (Rust), cmdargs (Haskell)

# General advice

- use *unit tests*
- use a *build system*, e.g. cargo (Rust), stack or cabal (Haskell)
- use a *cli library*, e.g. clap (Rust), cmdargs (Haskell)
- use a *parser combinator library*, e.g. nom (Rust), parsec (Haskell)

# General advice

- use *unit tests*
- use a *build system*, e.g. cargo (Rust), stack or cabal (Haskell)
- use a *cli library*, e.g. clap (Rust), cmdargs (Haskell)
- use a *parser combinator library*, e.g. nom (Rust), parsec (Haskell)
- use *iterators* to act lazily on large data (e.g. corpora)