

Tutorium 1: Grammar Induction

Praktikum Verarbeitung natürlicher Sprache

Richard Mörbitz

2021-04-19

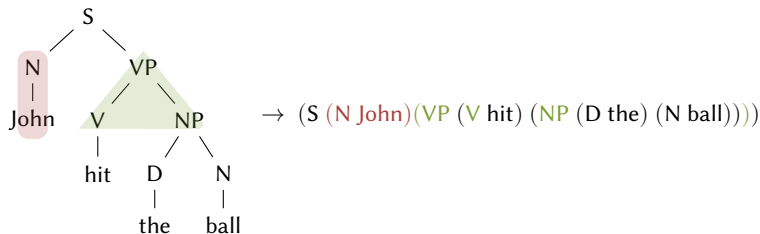
Overview

- 1 Introduction to the training corpus
- 2 Algorithmic approach
- 3 Overall pipeline (suggestion)
- 4 Data structures (suggestion)
- 5 General advice

Introduction to the training corpus

Take a look at `training.mrg`:

- one line = one constituent tree (annotated sentence)
- trees are stored as *symbolic expressions* (s-expressions)



Inductive definition of s-expressions:

atoms: strings that do not contain spaces or parenthesis

lists: $(s_1 \sqcup s_2 \sqcup \dots \sqcup s_k)$ where s_i are s-expressions (\sqcup is space)

Introduction to the training corpus: some remarks

Trees as s-expressions:

- no empty lists
- s_1 is always an atom

Parenthesis in grammar symbols are escaped (-LRB-, -RRB-).

Sentences are *tokenized*, e.g.,

“There’s no food”, he says.

becomes

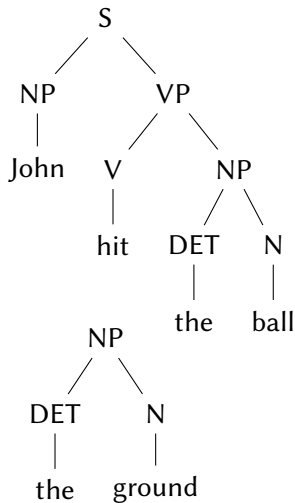
“ There ’s no food ” , he says .

Algorithmic approach

- 1: **function** INDUCE_GRAMMAR(corpus $c = t_1, \dots, t_n$)
- 2: **for** $i = 1, \dots, n$ **do**
- 3: **for each** position p of t_i **do**
- 4: $r \leftarrow$ readoff rule at p
- 5: increase count of r by 1
- 6: $R \leftarrow$ set of all counted rules
- 7: $N, \Sigma \leftarrow$ all nonterminals/terminals occurring in R
- 8: S depends on corpus (usually ROOT)
- 9: $p \leftarrow$ normalized counts of rules with same left-hand side
- 10: **return** $((N, \Sigma, S, R), p)$

Algorithmic approach

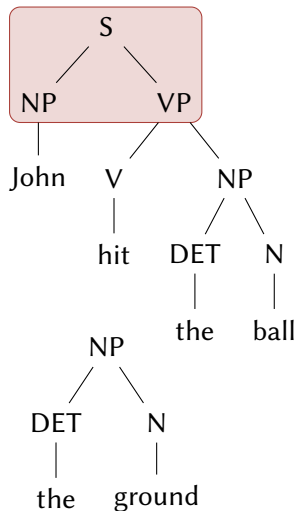
rules:



Algorithmic approach

rules:

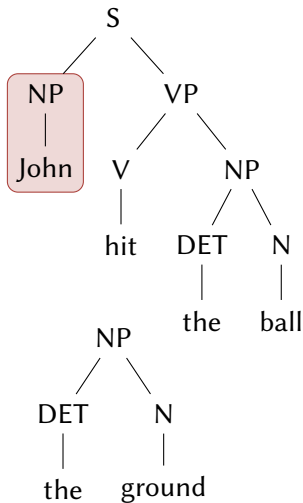
- $S \rightarrow NP VP$



Algorithmic approach

rules:

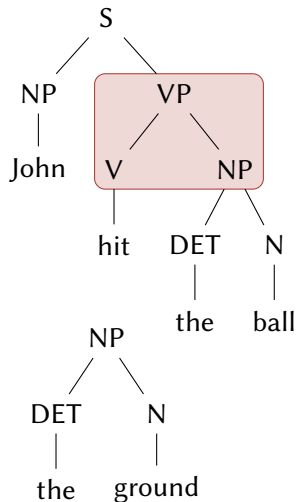
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$



Algorithmic approach

rules:

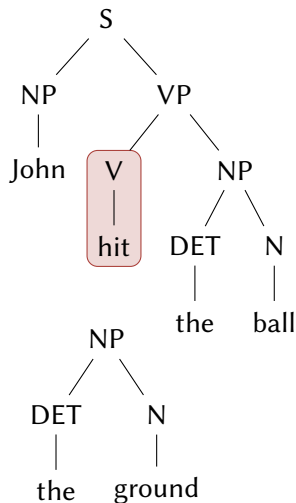
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$



Algorithmic approach

rules:

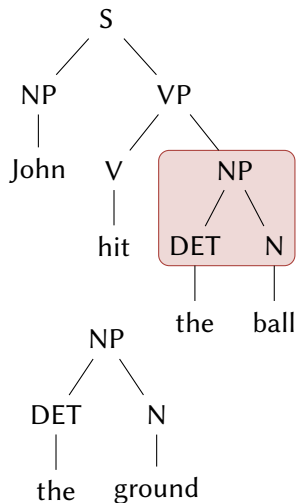
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$



Algorithmic approach

rules:

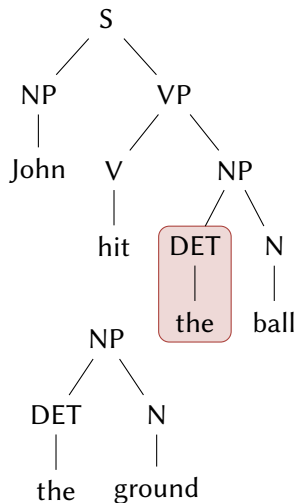
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$
- $NP \rightarrow \text{DET N}$



Algorithmic approach

rules:

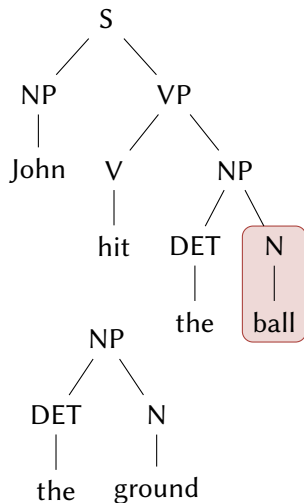
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$
- $NP \rightarrow \text{DET N}$
- $\text{DET} \rightarrow \text{the}$



Algorithmic approach

rules:

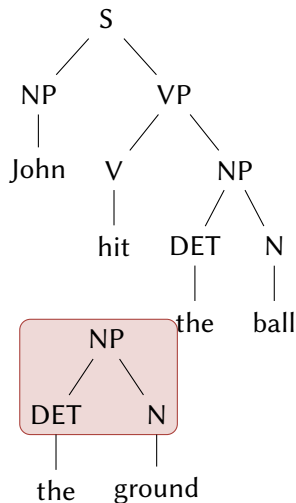
- $S \rightarrow NP VP$
- $NP \rightarrow \text{John}$
- $VP \rightarrow V NP$
- $V \rightarrow \text{hit}$
- $NP \rightarrow \text{DET N}$
- $\text{DET} \rightarrow \text{the}$
- $N \rightarrow \text{ball}$



Algorithmic approach

rules:

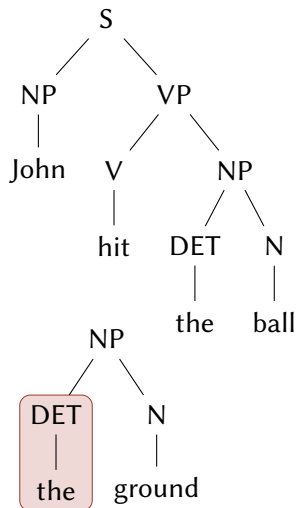
- $S \rightarrow NP VP$ # 1
- $NP \rightarrow John$ # 1
- $VP \rightarrow V NP$ # 1
- $V \rightarrow hit$ # 1
- $NP \rightarrow DET N$ # 2
- $DET \rightarrow the$ # 1
- $N \rightarrow ball$ # 1



Algorithmic approach

rules:

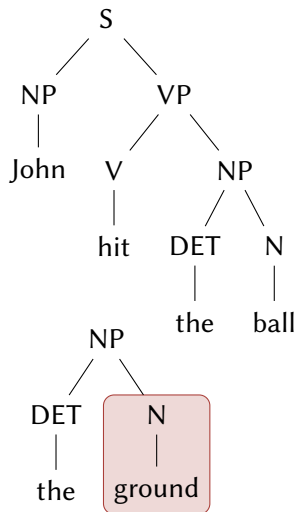
- $S \rightarrow NP VP$ # 1
- $NP \rightarrow John$ # 1
- $VP \rightarrow V NP$ # 1
- $V \rightarrow hit$ # 1
- $NP \rightarrow DET N$ # 2
- $DET \rightarrow the$ # 2
- $N \rightarrow ball$ # 1



Algorithmic approach

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



Algorithmic approach

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

Algorithmic approach

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
- $NP \rightarrow DET N$ # $\frac{2}{1+2} = 2/3$

Algorithmic approach

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
- $NP \rightarrow DET N$ # $\frac{2}{1+2} = 2/3$
- $DET \rightarrow the$ # $\frac{2}{2} = 1$

Algorithmic approach

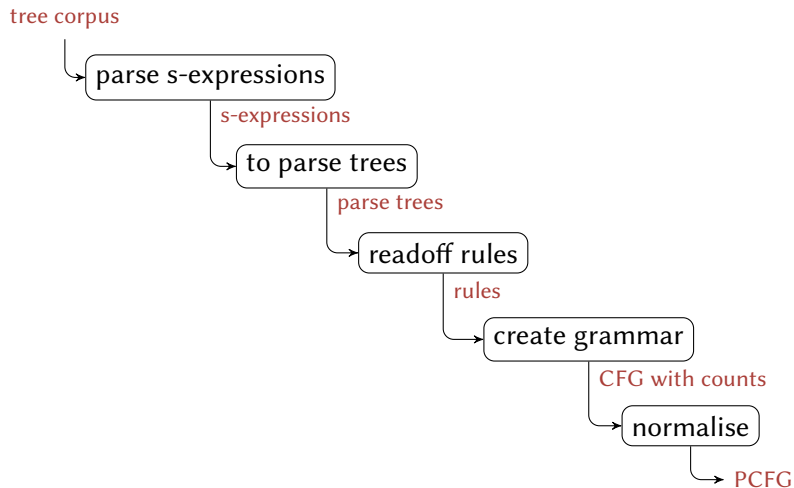
rules:

- $S \rightarrow NP VP$ # 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

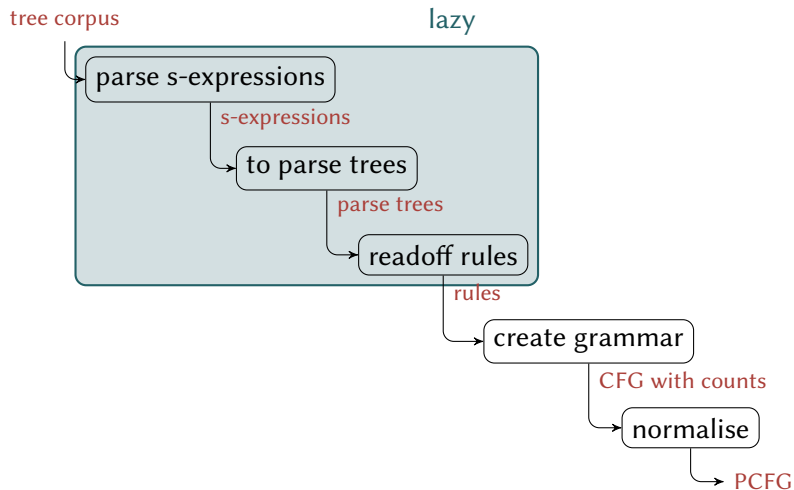
normalised rules:

- $S \rightarrow NP VP$ # 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)



Data structures (suggestion)

- s-expression:

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

A can be chosen the string type

Data structures (suggestion)

- s-expression:

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

- tree

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

A can be chosen the string type

Data structures (suggestion)

- rule

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

N and T can be chosen as the string type, w as a floating point type

Data structures (suggestion)

- 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 }
```

N and T can be chosen as the string type, W as a floating point type

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)