

# Zpracování přirozeného jazyka

**Jindřich Matuška**

Faculty of Informatics, Masaryk University

14. listopadu 2024

# Čas na odpovědníky

# Obsah

Předzpracování dat

Gramatiky

# Obsah

Předzpracování dat

Gramatiky

# Zdroje dat

Odkud můžeme čerpat data?

# Zdroje dat

Odkud můžeme čerpat data?

- Weby (HTML stránky – Odstranění značek, nebo z nich lze vytáhnout více dat?)
- Knihy (Získání textu ze stránek. Lze vytáhnout více dat?)
- Sociální sítě (Můžeme? Anonymizace?)

# Předzpracování dat

- Co je naším cílem?
- Co je pro náš cíl důležité, podstatné?
- Co z dat dokážeme vyčíst? Dokážeme vyčíst něco více než holý text?
- Jsou data dostatečně čistá? Je třeba je vyčistit? Lze je vyčistit?
- Anotace dat?
- V jaké formě budeme data vůbec ukládat?

# Obsah

Předzpracování dat

Gramatiky



# Bezkontextové gramatiky

- Množina terminálních symbolů  $\Sigma = \{a, b, c, \dots\}$
- Množina neterminálních symbolů  $N = \{A, B, C, \dots\}$
- Speciální neterminál (kořen)  $S \in N$
- Soubor pravidel  $\subseteq N \times V^*$ , kde  $V = \Sigma \cup N$ 
  - neterminál  $\rightarrow$  libovolný řetězec
  - Pokud neterminál uvozuje více pravidel, používáme svislítko |

# Příklad

$$S \rightarrow NP VP,$$
$$NP \rightarrow Noun \mid Ad NP,$$
$$VP \rightarrow Verb,$$
$$Noun \rightarrow \text{dítě} \mid \text{člověk} \mid \text{kapsa},$$
$$Adj \rightarrow \text{starý} \mid \text{cestující} \mid \text{nové},$$
$$Verb \rightarrow \text{píše} \mid \text{sedí} \mid \text{mluví},$$

# Syntaktický strom

- Kořenem je kořen gramatiky  $S$
- Listy jsou terminály
- Potomci každého uzlu jsou seřazení
- Pro každý uzel:
  - uzel je terminál, nebo
  - potomci uzlu jsou pravidlem gramatiky
- Může existovat více různých odvození

Tvorba syntaktického stromu z věty se nazývá syntaktický analýza

# Příklad

$$S \rightarrow NP VP,$$
$$NP \rightarrow Noun \mid Adj NP,$$
$$VP \rightarrow Verb,$$
$$Noun \rightarrow \text{dítě} \mid \text{člověk} \mid \text{kapsa},$$
$$Adj \rightarrow \text{starý} \mid \text{cestující} \mid \text{nové},$$
$$Verb \rightarrow \text{píše} \mid \text{sedí} \mid \text{mluví},$$

1. „cestující sedí“
2. „nové nové kapsa píše“
3. „starý člověk mluví“

## Příklad 9.2.1

Uvažte gramatiku s následujícími pravidly:

$$\begin{aligned} S &\rightarrow NP VP \mid VP, \\ NP &\rightarrow Noun \mid NP Conj NP, \\ VP &\rightarrow NP Esse \end{aligned}$$

a následujícím lexikonem:

$$\begin{aligned} Noun &\rightarrow \text{Romulus} \mid \text{Remus} \mid \text{Danubius} \mid \text{fratellus} \mid \text{fratelli} \mid \text{fluvius}, \\ Esse &\rightarrow \text{sum} \mid \text{est} \mid \text{sunt} \mid \text{eram} \mid \text{erat} \mid \text{erant}, \\ Conj &\rightarrow \text{et} \end{aligned}$$

a) Rozhodněte, které z následujících vět lze v gramatice vygenerovat.

1. „Romulus et Remus fratelli erant“
2. „Remus et Danubius et Romulus sum“
3. „Danubius est fluvius“

b) Nalezněte ke každé větě její syntaktický strom, pokud existuje.

# Pokrytí vs. přesnost

Pro zamýšlený jazyk  $L$  a gramatiku  $G$  generující jazyk  $L(G)$ :

- Pokrytí  $\frac{|L \cup L(G)|}{|L|}$
- Přesnost  $\frac{|L \cap L(G)|}{|L(G)|}$

## Příklad 9.2.3

Uvažme gramatiku  $G$  s následujícími pravidly

$$S \rightarrow AAA \mid BA \mid AB \mid C,$$

$$A \rightarrow a,$$

$$B \rightarrow bA \mid Ab,$$

$$C \rightarrow BA \mid AB \mid cB \mid Bc$$

a jazyk  $L$  vět délky 3 sestavených ze slov  $a, b, c$ , které obsahují slovo  $b$  právě jednou.

- a) Je gramatika jednoznačná, neboli existuje pro každou větu jazyka  $L(G)$  odvoditelnou v  $G$  právě jeden syntaktický strom?
- b) Jaké je pokrytí gramatiky  $G$  vzhledem k zamýšlenému jazyku  $L$ ?
- c) Jaká je přesnost gramatiky  $G$  vzhledem k zamýšlenému jazyku  $L$ ?