

Gramatická indukcia a jej využitie

Michal Malý

KAI FMFI UK

29. Marec 2010

Prehľad

- 1 Teória formálnych jazykov
- 2 Gramatická indukcia
- 3 Príklady gramatickej indukcie

Na počiatku bolo slovo. A slovo...

Definícia (Slovo)

Slovo je konečná postupnosť symbolov z abecedy.

Definícia (Abeceda)

Abeceda – množina znakov (symbolov, tokenov)

Definícia (Jazyk)

Jazyk – množina slov.

Príklad:

abeceda: $\Sigma = \{a, b\}$

slová: *abba,aaaabbbb*

jazyk: $L = \{\varepsilon, ab, ba, aabb, abba, baba, bbaa, \dots\}$ = tie slová, kde počet *a*-čok a *b*-čok v slove je rovnaký

Gramatika

Definícia (Terminálny symbol)

Terminály – symboly, z ktorých sa skladajú slová (na výstupe).

Definícia (neterminálny symbol)

Neterminály – pomocné symboly pri odvodzovaní, nesmú sa objaviť na výstupe.

Definícia (Gramatika)

Gramatika je určená množinou neterminálov, množinou terminálov, množinou odvodzovacích pravidiel a štartovacím symbolom.

Gramatika – príklad

terminály: $T = \{1, 2, 3, \dots, 9, +, -, (,)\}$

neterminály: $N = \{V, C\}$

štartovací symbol: V

pravidlá:

$V \rightarrow C$

$V \rightarrow V + V$

$V \rightarrow V - V$

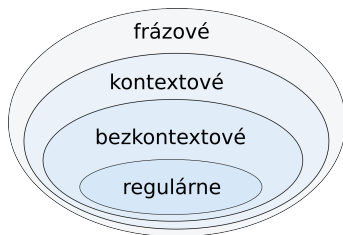
$V \rightarrow (V)$

$C \rightarrow 1, C \rightarrow 2, \dots, C \rightarrow 9$

Vieme vyrobiť slovo $1 - (4 + 5)$ napríklad takto:

$V \Rightarrow V - V \Rightarrow C - V \Rightarrow C - (V) \Rightarrow C - (V + V) \Rightarrow C - (C + V) \Rightarrow$
 $C - (C + C) \Rightarrow 1 - (C + C) \Rightarrow 1 - (4 + C) \Rightarrow 1 - (4 + 5)$

Chomského hierarchia



- 1 regulárne gramatiky \equiv konečný automat (nevedia $a^n b^n$)
- 2 bezkontextové gramatiky \equiv nedeterministický zásobníkový automat (nevedia $a^n b^n c^n$)
- 3 kontextové gramatiky \equiv nedeterministický lineárny automat (nevie EXPSPACE-hard problémy)
- 4 frázové gramatiky (neobmedzené, typu 0) \equiv Turingov stroj (vie všetko, čo je „intuitívne vypočítateľné“)

Krátka charakterizácia gramatík

- regulárna gramatika: vľavo iba neterminál, vpravo sa neterminál môže vyskytovať iba na konci: $S \rightarrow aX$
- bezkontextová gramatika: vľavo iba neterminál: $S \rightarrow XaYbb$
- kontextové gramatiky: môže sa vyskytovať kontext, ktorý ale ostáva: $\alpha S\beta \rightarrow \alpha XYa\beta$
- frázové gramatiky: neobmedzené

Gramatická indukcia

- Jazyk – zápis (formalizácia) prakticky ľubovoľného problému
- Príklady: pozitívne, negatívne
- Riešenie problému?

Definícia

Gramatická indukcia je spôsob, ako odvodiť formálnu gramatiku z množiny vzoriek – pozorovaní.

- využitie – spracovanie prirodzeného jazyka, kompresia, ...

Gramatická indukcia – metódy

- pokus-omyl
- greedy (LZW,Sequitur)
- genetické algoritmy (John R. Koza)
- formálne metódy

Gramatická indukcia – pokus-omyl

pozitívne príklady: $\mathcal{D}^+ = \{a, aaa, aaab, aab\}$

negatívne príklady: $\mathcal{D}^- = \{ab, abc, abb, aabb\}$

i	x_i^+	\mathcal{P}	\mathcal{P} produces \mathcal{D}^- ?
1	a	$S \rightarrow A$ $A \rightarrow a$	No
2	aaa	$S \rightarrow A$ $A \rightarrow a$ $A \rightarrow aA$	No
3	aaab	$S \rightarrow A$ $A \rightarrow a$ $A \rightarrow aA$ $A \rightarrow ab$	Yes: $ab \in \mathcal{D}^-$
3	aaab	$S \rightarrow A$ $A \rightarrow a$ $A \rightarrow aA$ $A \rightarrow aab$	No
4	aab	$S \rightarrow A$ $A \rightarrow a$ $A \rightarrow aA$ $A \rightarrow aab$	No

Gramatická indukcia – stemming cez Myhill-Nerodovu ekvivalenciu

Veta (Myhill-Nerodova veta)

Jazyk L je regulárny práve vtedy, ak relácia R definovaná ako

$$u R v \iff_{def} \forall x \in \Sigma^* (ux \in L \iff vx \in L)$$

je reláciou ekvivalencie konečného indexu.

- dve slová sú v relácii, ak pridaním nejakého suffixu (koncevky) dostanem pre obe buď slovo z jazyka, alebo pre obe slovo nie z jazyka
- konštrukcia minimálneho deterministického automatu

Výsledná gramatika pre slovo mesto

pravidlá pre neterminál číslo...	použití
0 → ... mesta 56 meste mesteck 61 mesto 48 mestsk 455 mestu ...	(štart)
3 → i ε	111
15 → ho j ε	179
48 → m ε	101
49 → ch m ε	111
56 → ch m 3 ε	121
61 → a o u	7
83 → m u	35
455 → a e 15 i o 83 u y 49	2

Výsledný stemovací slovník pre slovo mesto

...

mestach mestami mestam mesta

meste

mestecka mestecko mestecku

mestom mesto

mestska mestskeho mestskej mestske mestski mestskom mestskou

mestsku mestskych mestskym mestsky

...

Koniec

Ďakujem za pozornosť.
Otázky?