

Formální jazyky a automaty

Syntaktická analýza, přímá a nepřímá levá rekurze, Greibachové normální forma

Jan Křetínský

Fakulta informatiky, MU Brno

Jaro 2024

Bezkontextové gramatiky → zásobníkové automaty

Motivace 1: Jaká je třída jazyků rozpoznávaných zásobníkovými automaty?

Motivace 2: Je dána bezkontextová gramatika \mathcal{G} a slovo w .
Jak zjistit, zda se slovo w dá vygenerovat v gramatice \mathcal{G} ?

Problém syntaktické analýzy pro bezkontextové gramatiky:
pro danou bezkontextovou gramatiku \mathcal{G} a slovo w rozhodnout, zda $w \in L(\mathcal{G})$.

Věta 3.47.

Ke každé CFG \mathcal{G} lze sestrojít PDA \mathcal{M} takový, že $L(\mathcal{G}) = L_e(\mathcal{M})$.

Důkaz. K dané gramatice \mathcal{G} konstruujeme PDA \mathcal{M} simulující derivace v \mathcal{G}

Nedeterministická syntaktická analýza shora dolů

Věta 3.47.

Ke každé CFG \mathcal{G} lze sestavit PDA \mathcal{M} takový, že $L(\mathcal{G}) = L_e(\mathcal{M})$.

Důkaz. K dané gramatice \mathcal{G} konstruujeme PDA \mathcal{M} simulující levé derivace v \mathcal{G} , tj. vždy rozvíjíme nejlevější neterminál:

- ▶ V levé derivaci je v jednom kroku odvození nahrazen (nejlevější) neterminál A pravou stranou $X_1 \dots X_n$ nějakého A -pravidla.
- ▶ V \mathcal{M} této situaci odpovídá náhrada A na vrcholu zásobníku řetězem $X_1 \dots X_n$.

$\mathcal{M} = (\{q\}, \Sigma, N \cup \Sigma, \delta, q, S, \emptyset)$, kde δ je definována:

- ▶ (expanze) $\delta(q, \varepsilon, A)$ obsahuje (q, α) právě když $A \rightarrow \alpha \in P$
- ▶ (čtení) $\delta(q, a, a) = \{(q, \varepsilon)\}$ pro všechna $a \in \Sigma$

Příklad

$S \rightarrow aAB$

$A \rightarrow Aa \mid \varepsilon$

$B \rightarrow SaA \mid b$

$\delta(q, \varepsilon, S) = \{(q, aAB)\}$

$\delta(q, \varepsilon, A) = \{(q, Aa), (q, \varepsilon)\}$

$\delta(q, \varepsilon, B) = \{(q, SaA), (q, b)\}$

$\delta(q, a, a) = \delta(q, b, b) = \{(q, \varepsilon)\}$

$S \Rightarrow aAB$

$\Rightarrow aB$

$\Rightarrow aSaA$

$\Rightarrow aaABaA$

$\Rightarrow aaBaA$

$\Rightarrow aabaA$

$\Rightarrow aaba$

$$A \Rightarrow^* w \iff (q, w, A) \vdash^* (q, \varepsilon, \varepsilon)$$

(\implies) Indukcí vzhledem k délce odvození m .

- ▶ $m = 0$: zřejmé.
- ▶ $m \geq 1$: necht' tvrzení platí pro všechna $m' < m$.

$A \Rightarrow X_1 X_2 \dots X_k \Rightarrow^* x_1 x_2 \dots x_k = w$, kde $X_i \xRightarrow{m_i} x_i$, $0 \leq m_i < m$
z definice δ plyne $(q, w, A) \vdash (q, w, X_1 X_2 \dots X_k)$.

Je-li $X_i \in N$, pak dle indukčního předpokladu máme
 $(q, x_i, X_i) \vdash^* (q, \varepsilon, \varepsilon)$.

Je-li $X_i \in \Sigma$, pak $X_i = x_i$ a z definice δ plyne
 $(q, x_i, x_i) \vdash (q, \varepsilon, \varepsilon)$.

Kompozicí dostáváme $(q, w, A) \vdash^+ (q, \varepsilon, \varepsilon)$.

(\Leftarrow) Předpokládejme $(q, w, A) \vdash^n (q, \varepsilon, \varepsilon)$ a ukažme $A \Rightarrow^* w$.
Indukcí vzhledem k délce výpočtu n .

▶ $n = 0$: zřejmé.

▶ $n \geq 1$: nechť tvrzení platí pro všechna $n' < n$.

$(q, w, A) \vdash (q, w, X_1 X_2 \dots X_k)$, tj. $A \rightarrow X_1 X_2 \dots X_k \in P$

w můžeme napsat jako $w = x_1 x_2 \dots x_k$ takové, že

▶ je-li $X_i \in N$, pak $(q, x_i, X_i) \vdash^{n_i} (q, \varepsilon, \varepsilon)$, kde $n_i < n$.

Dle IP $X_i \Rightarrow^+ x_i$.

▶ je-li $X_i \in \Sigma$, pak $X_i \xrightarrow{0} x_i$.

Vhodnou kompozicí obdržíme

$A \Rightarrow X_1 X_2 \dots X_k \Rightarrow^* x_1 X_2 \dots X_k \Rightarrow^* x_1 \dots x_k = w$

což je levá derivace slova w v gramatice \mathcal{G} . □

Nedeterministická syntaktická analýza zdola nahoru

$S \rightarrow XY$

$X \rightarrow ab$

$Y \rightarrow c$

Věta 3.55.

Nechť \mathcal{G} je libovolná CFG, pak lze zkonstruovat rozšířený PDA \mathcal{R} takový, že $L(\mathcal{G}) = L(\mathcal{R})$.

Důkaz. Vrchol zásobníku píšeme vpravo.

Konstruujeme rozšířený PDA \mathcal{R} , který simuluje **pravou** derivaci v \mathcal{G} **v obráceném pořadí**.

PDA \mathcal{R} má kroky dvojího typu:

1. (**čtení**) může kdykoli načíst do zásobníku symbol ze vstupu
2. (**redukce**) je-li na vrcholu zásobníku řetězec tvořící pravou stranu nějakého pravidla v \mathcal{G} , může ho nahradit odpovídajícím levostranným neterminálem (a ze vstupu nic nečte)

Nedeterministická syntaktická analýza zdola nahoru

Věta 3.55.

Nechť \mathcal{G} je libovolná CFG, pak lze zkonstruovat rozšířený PDA \mathcal{R} takový, že $L(\mathcal{G}) = L(\mathcal{R})$.

Důkaz. Vrchol zásobníku píšeme vpravo.

Konstruujeme rozšířený PDA \mathcal{R} , který simuluje **pravou** derivaci v \mathcal{G} **v obráceném pořadí**.

PDA \mathcal{R} má kroky dvojího typu:

1. (**čtení**) může kdykoli načíst do zásobníku symbol ze vstupu
2. (**redukce**) je-li na vrcholu zásobníku řetězec tvořící pravou stranu nějakého pravidla v \mathcal{G} , může ho nahradit odpovídajícím levostranným neterminálem (a ze vstupu nic nečte)

Nechť $\mathcal{G} = (N, \Sigma, P, S)$.

Položme $\mathcal{R} = (\{q, r\}, \Sigma, N \cup \Sigma \cup \{\perp\}, \delta, q, \perp, \{r\})$, kde \perp je nově přidaný symbol a kde δ je definována takto:

1. $\delta(q, a, \varepsilon) = \{(q, a)\}$ pro všechna $a \in \Sigma$
2. je-li $A \rightarrow \alpha$ pravidlo v P , pak $\delta(q, \varepsilon, \alpha)$ obsahuje (q, A)
3. $\delta(q, \varepsilon, \perp S) = \{(r, \varepsilon)\}$

Příklad

	krok výpočtu	odpovídající pravidlo z \mathcal{G}
$(q, i + i * i, \perp)$	$\vdash^i (q, +i * i, \perp i)$	$F \rightarrow i$
	$\vDash^{\varepsilon} (q, +i * i, \perp F)$	$T \rightarrow F$
	$\vDash^{\varepsilon} (q, +i * i, \perp T)$	$E \rightarrow T$
	$\vDash^{\varepsilon} (q, +i * i, \perp E)$	
	$\vDash^+ (q, \quad i * i, \perp E+)$	
	$\vdash^i (q, \quad *i, \perp E + i)$	$F \rightarrow i$
	$\vDash^{\varepsilon} (q, \quad *i, \perp E + F)$	$T \rightarrow F$
	$\vDash^{\varepsilon} (q, \quad *i, \perp E + T)$	
	$\vDash^* (q, \quad i, \perp E + T^*)$	
	$\vdash^i (q, \quad \varepsilon, \perp E + T * i)$	$F \rightarrow i$
	$\vDash^{\varepsilon} (q, \quad \varepsilon, \perp E + T * F)$	$T \rightarrow T * F$
	$\vDash^{\varepsilon} (q, \quad \varepsilon, \perp E + T)$	$E \rightarrow E + T$
	$\vDash^{\varepsilon} (q, \quad \varepsilon, \perp E)$	
	$\vDash^{\varepsilon} (r, \quad \varepsilon, \varepsilon)$	

$$S \Rightarrow^* \alpha A y \xRightarrow{n} xy \iff (q, xy, \perp) \vdash^* (q, y, \perp \alpha A),$$

kde $S \Rightarrow^* \alpha A y \xRightarrow{n} xy$ je pravá derivace a A je nejpravější neterminál.

(\implies) indukcí k délce odvození

(\impliedby) indukcí k délce výpočtu

Pro $A = S$ a $\alpha, y = \varepsilon$ dostáváme:

$$S \Rightarrow^* x \iff (q, x, \perp) \vdash^* (q, \varepsilon, \perp S) \quad \left[\vdash (r, \varepsilon) \right]$$



Efektivnost syntaktické analýzy

Nederministický PDA \implies nedeterministický algoritmus
 \implies exponenciální deterministický algoritmus

Řešení:

- ▶ deterministický algoritmus složitosti $\mathcal{O}(n^3)$, kde $n = |w|$ (algoritmus Cocke - Younger - Kasami)
- ▶ deterministické zásobníkové automaty a deterministické bezkontextové jazyky
- ▶ lineární algoritmy pro speciální třídy deterministických bezkontextových jazyků

Kanonické tvary bezkontextových gramatik

- ▶ redukované bezkontextové gramatiky
- ▶ gramatiky bez ϵ -pravidel
- ▶ gramatiky bez jednoduchých pravidel
- ▶ vlastní gramatiky
- ▶ Chomského normální forma
- ▶ gramatiky bez levé rekurze
- ▶ Greibachové normální forma

Definice 3.28.

Neterminál A v CFG $\mathcal{G} = (N, \Sigma, P, S)$ se nazývá **levorekurzivní** jestliže v \mathcal{G} existuje derivace $A \Rightarrow^+ A\beta$.

CFG bez levorekurzivních neterminálů se nazývá **nelevorekurzivní**.

Je-li v CFG pravidlo tvaru $A \rightarrow A\alpha$, hovoříme o **přímé levé rekurzi** na A .

Praktický význam: některé nástroje pro automatickou tvorbu parserů k zadaným gramatikám vyžadují na vstupu nelevorekurzivní gramatiku (např. ANTLR).

Algoritmus odstranění přímé levé rekurze

Nechť CFG $\mathcal{G} = (N, \Sigma, P, S)$ je necyklická a bez ε -pravidel, v níž všechna A -pravidla (pravidla mající na levé straně A) jsou tvaru

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n,$$

kde každý řetěz β_i začíná symbolem různým od A .

Algoritmus odstranění přímé levé rekurze

Nechť CFG $\mathcal{G} = (N, \Sigma, P, S)$ je necyklická a bez ε -pravidel, v níž všechna A -pravidla (pravidla mající na levé straně A) jsou tvaru

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n,$$

kde každý řetěz β_i začíná symbolem různým od A .

Nechť $\mathcal{G}' = (N \cup \{A'\}, \Sigma, P', S)$, kde P' obdržíme z P tak, že všechna výše uvedená pravidla nahradíme pravidly:

$$A \rightarrow \beta_1 \mid \dots \mid \beta_n \mid \beta_1 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 \mid \dots \mid \alpha_m \mid \alpha_1 A' \mid \dots \mid \alpha_m A'$$

Algoritmus odstranění přímé levé rekurze

Nechť CFG $\mathcal{G} = (N, \Sigma, P, S)$ je necyklická a bez ε -pravidel, v níž všechna A -pravidla (pravidla mající na levé straně A) jsou tvaru

$$A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n,$$

kde každý řetěz β_i začíná symbolem různým od A .

Nechť $\mathcal{G}' = (N \cup \{A'\}, \Sigma, P', S)$, kde P' obdržíme z P tak, že všechna výše uvedená pravidla nahradíme pravidly:

$$A \rightarrow \beta_1 \mid \dots \mid \beta_n \mid \beta_1 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 \mid \dots \mid \alpha_m \mid \alpha_1 A' \mid \dots \mid \alpha_m A'$$

Pak $L(\mathcal{G}) = L(\mathcal{G}')$ a \mathcal{G}' je necyklická a bez ε -pravidel.

Lemma o substituci

Lemma 3.20. (o substituci)

Nechť $\mathcal{G} = (N, \Sigma, P, S)$ je CFG. Nechť $A \rightarrow \alpha_1 B \alpha_2 \in P$.

Nechť $B \rightarrow \beta_1 \mid \dots \mid \beta_r$ jsou všechna pravidla v P tvaru $B \rightarrow \alpha$.

Definujme $\mathcal{G}' = (N, \Sigma, P', S)$, kde

$$P' = (P \setminus \{A \rightarrow \alpha_1 B \alpha_2\}) \cup \{A \rightarrow \alpha_1 \beta_1 \alpha_2 \mid \dots \mid \alpha_1 \beta_r \alpha_2\}.$$

Pak $L(\mathcal{G}) = L(\mathcal{G}')$.

Lemma o substituci

Lemma 3.20. (o substituci)

Nechť $\mathcal{G} = (N, \Sigma, P, S)$ je CFG. Nechť $A \rightarrow \alpha_1 B \alpha_2 \in P$.

Nechť $B \rightarrow \beta_1 \mid \dots \mid \beta_r$ jsou všechna pravidla v P tvaru $B \rightarrow \alpha$.

Definujme $\mathcal{G}' = (N, \Sigma, P', S)$, kde

$$P' = (P \setminus \{A \rightarrow \alpha_1 B \alpha_2\}) \cup \{A \rightarrow \alpha_1 \beta_1 \alpha_2 \mid \dots \mid \alpha_1 \beta_r \alpha_2\}.$$

Pak $L(\mathcal{G}) = L(\mathcal{G}')$.

Příklad:

$A \rightarrow Bd \mid c$

$B \rightarrow Bdd \mid Ccc \mid aAd$

$C \rightarrow Aa$

Algoritmus odstranění levé rekurze

Vstup: Vlastní CFG $\mathcal{G} = (N, \Sigma, P, S)$

Výstup: Ekvivalentní nelevorekurzivní gramatika bez ε -pravidel

- 1: Uspořádej libovolně N , $N = \{A_1, \dots, A_n\}$
- 2: **for** $i \leftarrow 1$ **to** n **do**
- 3: **for** $j \leftarrow 1$ **to** $i - 1$ **do**
- 4: **for all** pravidlo tvaru $A_i \rightarrow A_j\alpha$ **do**
- 5: přidej pravidla $A_i \rightarrow \beta_1\alpha \mid \dots \mid \beta_k\alpha$
- 6: (kde $A_j \rightarrow \beta_1 \mid \dots \mid \beta_k$ jsou všechna pravidla pro A_j)
- 7: vypusť pravidlo $A_i \rightarrow A_j\alpha$
- 8: odstraň případnou přímou levou rekurzi na A_i

Algoritmus odstranění levé rekurze

Vstup: Vlastní CFG $\mathcal{G} = (N, \Sigma, P, S)$

Výstup: Ekvivalentní nelevorekurzivní gramatika bez ε -pravidel

- 1: Uspořádej libovolně N , $N = \{A_1, \dots, A_n\}$
- 2: **for** $i \leftarrow 1$ **to** n **do**
- 3: **for** $j \leftarrow 1$ **to** $i - 1$ **do**
- 4: **for all** pravidlo tvaru $A_i \rightarrow A_j\alpha$ **do**
- 5: přidej pravidla $A_i \rightarrow \beta_1\alpha \mid \dots \mid \beta_k\alpha$
- 6: (kde $A_j \rightarrow \beta_1 \mid \dots \mid \beta_k$ jsou všechna pravidla pro A_j)
- 7: vypuť pravidlo $A_i \rightarrow A_j\alpha$
- 8: odstraň případnou přímou levou rekurzi na A_i

Příklad:

$A \rightarrow Bd \mid c$

$B \rightarrow Bdd \mid Ccc \mid aAd$

$C \rightarrow Aa$

Korektnost algoritmu

Konečnost.

Ekvivalence gramatik: Všechny úpravy jsou dle Lemmatu o substituci nebo odstraňují přímou levou rekurzi.

Výsledná gramatika je nelevorekurzivní:

1. po i -té iteraci vnějšího cyklu začíná každé A_i -pravidlo buď terminálem nebo neterminálem A_k , kde $k > i$.
2. po j -té iteraci vnitřního cyklu začíná každé A_i -pravidlo buď terminálem nebo neterminálem A_k , kde $k > j$.

Výsledná gramatika je bez ε -pravidel.

Greibachové normální forma

Definice 3.33.

Bezkontextová gramatika $\mathcal{G} = (N, \Sigma, P, S)$ je v **Greibachové normální formě** (GNF), právě když

- ▶ \mathcal{G} je bez ε -pravidel a
- ▶ každé pravidlo z P je tvaru $A \rightarrow a\alpha$, kde $a \in \Sigma$ a $\alpha \in N^*$ (s případnou výjimkou pravidla $S \rightarrow \varepsilon$).

Věta 3.34.

Každý bezkontextový jazyk lze generovat bezkontextovou gramatikou v Greibachové normální formě.

Důkaz.

$L = \emptyset$: zřejmé ($S \rightarrow aS$)

- $L \neq \emptyset$:
1. z vlastní gramatiky eliminujeme levou rekurzi
 2. pak převedeme do GNF



Příklad

$A \rightarrow Ba \mid Db \mid c$

$B \rightarrow CC$

$C \rightarrow aE$

$D \rightarrow CDa \mid Eb$

$E \rightarrow bb$

Algoritmus transformace do GNF

Vstup: Nelevorekurzivní CFG $\mathcal{G} = (N, \Sigma, P, S)$ bez ε -pravidel

Výstup: Ekvivalentní gramatika v GNF

- 1: Najdi lineární uspořádání \prec splňující $(A \rightarrow B\alpha) \in P \implies A \prec B$
- 2: Označme $N = \{A_1, \dots, A_n \mid A_{i-1} \prec A_i, 1 < i \leq n\}$
- 3: **for** $i \leftarrow n - 1$ **downto** 1 **do**
- 4: **for all** pravidlo tvaru $A_i \rightarrow A_j\alpha$, kde $j > i$ **do**
- 5: přidej pravidlo $A_i \rightarrow \beta_1\alpha \mid \dots \mid \beta_k\alpha$
- 6: (kde $A_j \rightarrow \beta_1 \mid \dots \mid \beta_k$ jsou všechna A_j -pravidla)
- 7: vypuť pravidlo $A_i \rightarrow A_j\alpha$
- 8: Nahraď potřebné terminály novými neterminály a přidej příslušná pravidla

Korektnost algoritmu

Konečnost.

Ekvivalence gramatik.

Výsledná gramatika je v GNF.