

Bezkontextové gramatiky, derivačné stromy, úpravy gramatík

Ing. Viliam Hromada, PhD.

C-510
Ústav informatiky a matematiky
FEI STU

`viliam.hromada@stuba.sk`



Pre pripomenutie:

- **Regulárne jazyky** - jazyky generovateľné pomocou tzv. regulárnych gramatík.
- **Bezkontextové jazyky** - jazyky generovateľné pomocou tzv. bezkontextových gramatík.



Jednoznačné regexy

Našťastie, v tomto prípade existuje aj jednoznačná gramatika pre regulárne výrazy. Uvedieme príklad gramatiky pre regexy, ktoré budú popisovať reťazce nad abecedou $\{a, b\}$. **POZOR!!! Všimnite si, že v tejto gramatike sú zvislá čiara $|$, \emptyset , hviezdička či ε terminálne symboly - čo je v protiklade s tým, ako sme ich doteraz používali!).** Je to dané tým, že chceme popísať gramatiku, ktorá generuje reťazce predstavujúce regulárne výrazy a v nich majú tieto symboly svoj špecifický význam.

Pre jednoduchosť, na rozdiel od klasickej definície regexov, naša gramatika nebude vedieť generovať regexy, ktoré obsahujú symbol prázdnej množiny.

Gramatika $G =$

$(\{\langle \text{Regex} \rangle, \langle \text{Term} \rangle, \langle \text{Factor} \rangle, \langle \text{Base} \rangle, \langle \text{Char} \rangle\}, \{*, |, a, b, \varepsilon, \emptyset, (,)\}, P, \langle \text{Regex} \rangle)$,

kde pravidlá sú na ďalšom slajde:



Ľavé (pravé) odvodenie

Definícia

*Nech G je bezkontextová gramatika. Potom odvodenie $\alpha_0 \Rightarrow \alpha_1 \Rightarrow \dots \Rightarrow \alpha_k$ v G , kde každý krok odvodenia $\alpha_i \Rightarrow \alpha_{i+1}$ pre $0 \leq i \leq k$ realizujeme tak, že v reťazci α_i nahradíme prvý neterminál zľava, nazývame **ľavé odvodenie** (**ľavá derivácia**); ak nahrádzame prvý neterminál sprava, nazývame **pravé odvodenie** (**pravá derivácia**). Príslušné vznikajúce vetné formy nazývame **ľavé (pravé) vetné formy**. Krok odvodenia v ľavom (pravom) odvodení budeme označovať symbolom \Rightarrow_l (\Rightarrow_r).*



Transformácie gramatík

- Gramatiky možno **transformovať** na **ekvivalentné** gramatiky (t.j. generujúce ten istý jazyk), ktoré budú navyše **spĺňať** dodatočné **požiadavky** (napr. na tvar pravidiel).
- Základom pri transformácii gramatík je, aby sa po transformácii **nezmenil jazyk generovaný gramatikou**.



Odstránenie nadbytočných symbolov

V gramatike sa môžu vyskytovať symboly (N, T) ktoré sa alebo nikdy počas derivácie nepoužijú, alebo ak sa použijú, nikdy sa nepodarí odvodiť reťazec terminálov. Ich prítomnosť v gramatike potom zbytočne zväčšuje jej veľkosť a navyše môžu viesť k zlým výsledkom rôznych algoritmov, pracujúcich s gramatikami.

Definícia

Nech $G = (N, T, P, S)$ je bezkontextová gramatika. Symbol $X \in N \cup T$ sa nazýva:

- **nedostupný**, ak v G neexistuje odvodenie $S \Rightarrow^* \alpha X \beta$ pre nejaké $\alpha, \beta \in (N \cup T)^*$,
- **nadbytočný**, ak v G neexistuje odvodenie $S \Rightarrow^* xXz \Rightarrow^* xyz$ pre nejaké $x, y, z \in T^*$.



Odstránenie nadbytočných symbolov

Definícia

*Gramatika bez nadbytočných symbolov sa nazýva **redukovaná** gramatika.*

Každý nedostupný symbol je zároveň aj nadbytočný. Vytvorenie redukovanej gramatiky pozostáva z 2 hlavných fáz:

1. Odstránenie neterminálov, z ktorých nie je možné odvodiť reťazec terminálov alebo ϵ .
2. Odstránenie neterminálov a terminálov, ktoré sa nemôžu vyskytnúť vo vetných formách gramatiky.

Symbols sa odstránia spolu s príslušnými pravidlami, v ktorých vystupujú.



Odstránenie nadbytočných symbolov - množina N_T

V prvej fáze hľadáme neterminály, z ktorých je možné odvodiť reťazec terminálov (prípadne ε). Hľadáme množinu:

$$N_T = \{A \mid A \Rightarrow^* w, w \in T^*\} \quad (1)$$



Odstránenie nadbytočných symbolov - množina N_T

Vstup: bezkontextová gramatika $G = (N, T, P, S)$

Výstup: množina N_T podľa (1)

1: $N_T \leftarrow \emptyset$

2: **opakuj**

3: $N'_T \leftarrow N_T$

4: $N_T \leftarrow N'_T \cup \{A \mid A \rightarrow \alpha \in P \wedge \alpha \in (N'_T \cup T)^*\}$

5: **pokiaľ** $N_T \neq N'_T$



Odstránenie nadbytočných symbolov - množina N_T

Príklad: Nech $G = (\{S, A, B, C\}, \{a, b, c\}, P, S)$, kde P obsahuje pravidlá:

$$S \rightarrow ABc \mid abA$$

$$A \rightarrow aAa \mid bBb \mid b$$

$$B \rightarrow aBb \mid caBC$$

$$C \rightarrow \varepsilon \mid ab \mid BcCa$$

Podľa algoritmu by sa premenná N_T postupne menila: $\emptyset, \{A, C\}, \{S, A, C\}$.

Neterminál B je teda taký, že sa z neho nedá odvodiť terminálny reťazec. Jeho výskyt v gramatike je teda **nadbytočný**, pretože jeho odstránením sa **nijako nezmení jazyk** $L(G)$, ktorý gramatika generuje.



Odstránením neterminálu B a pravidiel s jeho výskytom by sme dostali gramatiku:

$$S \rightarrow abA$$

$$A \rightarrow aAa \mid b$$

$$C \rightarrow \varepsilon \mid ab$$



Odstránenie nadbytočných symbolov - množina V_D

V druhej fáze hľadáme neterminály a terminály, ktoré sa môžu vyskytnúť v nejakej vetnej forme danej gramatiky, t.j. sú to **dostupné symboly gramatiky**. Hľadáme množinu:

$$V_D = \{X \mid S \Rightarrow^* \alpha X \beta, \alpha, \beta \in (N \cup T)^*\}. \quad (2)$$



Odstránenie nadbytočných symbolov - množina V_D

Vstup: bezkontextová gramatika $G = (N, T, P, S)$

Výstup: množina V_D podľa (2)

1: $V_D \leftarrow \{S\}$

2: **opakuj**

3: $\acute{V}_D \leftarrow V_D$

4: $V_D \leftarrow \acute{V}_D \cup \{X \mid A \rightarrow \alpha X \beta \in P \wedge A \in \acute{V}_D\}$

5: **pokiaľ** $V_D \neq \acute{V}_D$



Odstránenie nadbytočných symbolov - množina V_D

Príklad: Nech $G = (\{S, A, B, C, D\}, \{a, b, c, d\}, P, S)$, kde P obsahuje pravidlá:

$$S \rightarrow AB \mid abS \mid \varepsilon$$

$$A \rightarrow aAa \mid b$$

$$B \rightarrow aCb \mid cC$$

$$C \rightarrow cS$$

$$D \rightarrow abC \mid \varepsilon \mid dD$$

Podľa algoritmu by sa premenná V_D postupne menila:

$\{S\}$, $\{S, A, B, a, b\}$, $\{S, A, B, C, a, b, c\}$. Nedostupné symboly sú teda D, d .



Odstránenie nadbytočných symbolov - výsledný algoritmus

Vstup: bezkontextová gramatika $G = (N, T, P, S)$

Výstup: gramatika bez nadbytočných symbolov

- 1: Vytvor množinu N_T .
- 2: **ak** $S \in N_T$ **potom**
- 3: odstráň z gramatiky G všetky neterminály, ktoré nepatria do N_T ;
- 4: vytvor množinu V_D
- 5: odstráň z gramatiky G všetky symboly, ktoré nepatria do V_D ;
- 6: **koniec ak**

Poradie krokov je **dôležité!**. Najprv sa musia odstrániť neterminály, ktoré nevedú na terminálne reťazce (N_T) a až potom sa odstránia nedostupné symboly (V_D).



Odstránenie nadbytočných symbolov

Príklad: Pokračujme v gramatike, pre ktorú sme hľadali množinu N_T na slajde č. 28. Po odstránení neterminálu B (pretože nebol v množine N_T) dostávame gramatiku $G_1 = (\{S, A, C\}, \{a, b, c\}, P_1, S)$, kde P_1 obsahuje pravidlá:

$$S \rightarrow abA$$

$$A \rightarrow aAa \mid b$$

$$C \rightarrow ab \mid \varepsilon$$

Pre túto gramatiku je množina $V_D = \{S, A, a, b\}$, t.j. symboly C, c sú nedostupné. Po ich odstránení dostávame redukovanú gramatiku:

$$G_{red} = (\{S, A\}, \{a, b\}, \{S \rightarrow abA, A \rightarrow aAa \mid b\}, S).$$

T.j. $L(G) = L(G_1) = L(G_{red}) = \{aba^nba^n \mid n \in \mathbb{Z}_0^+\}$.



Ak by sme gramatiku zo slajdu č. 28 upravovali **v opačnom poradí**, potom:

1. Množina $V_D = \{S, A, B, C, a, b, c\}$, t.j. všetky symboly sú dostupné a nič sa neodstráni.
2. Množina $N_T = \{A, C, S\}$, t.j. odstránime len neterminál B a výsledná gramatika:

$$S \rightarrow abA$$

$$A \rightarrow aAa \mid b$$

$$C \rightarrow \varepsilon \mid ab$$

pričom vidíme, že neterminál C nie je dostupný a terminál c ani len nevystupuje v niektorom z pravidiel.



Odstránenie ε -pravidiel

Definícia

Gramatika $G = (N, T, P, S)$ sa nazýva **gramatika bez ε -pravidiel**, ak množina pravidiel P neobsahuje žiadne ε -pravidlo (také, kde na pravej strane je len prázdny reťazec ε), alebo v P existuje jediné ε -pravidlo $S \rightarrow \varepsilon$, pričom začiatočný symbol gramatiky sa nevyskytuje na pravej strane žiadneho pravidla.



Odstránenie ε -pravidiel

- Každé ε -pravidlo odstránime a zároveň do gramatiky doplníme ďalšie pravidlá, aby sa nezmenil výsledný jazyk generovaný gramatikou.
- Pri dopĺňaní pravidiel postupujeme podľa toho, ktoré neterminály sa mohli počas derivácie prepísať na ε a teda sa mohli "*stratiť*" z vetnej formy.
- Jediný prípad, kedy akceptujeme ε -pravidlo je v prípade, že $\varepsilon \in L(G)$. V takom prípade doplníme nový počiatočný symbol S' , doplníme pravidlá $S' \rightarrow S \mid \varepsilon$ a ostatné ε -pravidlá príslušným spôsobom odstránime.



Príklad

Gramatiku s ε -pravidlami

$$S \rightarrow aSb \mid \varepsilon \mid A$$

$$A \rightarrow bAa \mid \varepsilon$$

Upravíme na:

$$S' \rightarrow S \mid \varepsilon$$

$$S \rightarrow aSb \mid A \mid ab$$

$$A \rightarrow bAa \mid ba$$



Odstránenie ε -pravidiel

Uvažujme o gramatike $G = (N, T, P, S)$. Hľadáme množinu neterminálov N_ε , ktoré sa môžu počas odvodzovania „stratiť“, t.j. prepísať na prázdne slovo ε .

$$N_\varepsilon = \{A \mid A \Rightarrow^* \varepsilon, A \in N\} \quad (3)$$

Množinu iteratívne tvoria tie neterminály, z ktorých možno odvodiť len ε .



Odstránenie ε -pravidiel

Vstup: bezkontextová gramatika $G = (N, T, P, S)$

Výstup: množina N_ε (3)

1: $N_\varepsilon \leftarrow \emptyset$.

2: **opakuj**

3: $N'_\varepsilon \leftarrow N_\varepsilon$

4: $N_\varepsilon \leftarrow N'_\varepsilon \cup \{A \mid A \rightarrow \alpha \in P \wedge \alpha \in N'^*_\varepsilon\}$

5: **pokiaľ** $N_\varepsilon \neq N'_\varepsilon$



Odstránenie ε -pravidiel

- Po určení N_ε je potrebné vytvoriť množinu pravidiel \hat{P} z množiny P tak, že sa v \hat{P} nebudú vyskytovať ε -pravidlá, ale nezmení sa generovaný jazyk.
- Do \hat{P} môžem dať tie pravidlá z P , ktoré na pravej strane neobsahujú neterminály z N_ε .
- Pre tie pravidlá z P , ktoré na pravej strane obsahujú neterminály z N_ε , t.j. tvaru:

$$A \rightarrow \alpha_0 B_1 \alpha_1 \dots B_k \alpha_k,$$

kde $\alpha_j \in ((N - N_\varepsilon) \cup T)^*$, $B_j \in N_\varepsilon$, platí, že každé takéto pravidlo nahradíme v \hat{P} množinou pravidiel tvaru

$$A \rightarrow \alpha_0 X_1 \alpha_1 \dots X_k \alpha_k,$$

kde $X_i = B_i$ alebo $X_i = \varepsilon$ pre $i = 1, \dots, k$.



Odstránenie ε -pravidiel

- POZOR! Ak by malo vzniknúť pravidlo $A \rightarrow \varepsilon$ alebo $A \rightarrow A$, tak ho do \hat{P} nepridáme.
- Navyše, ak $\varepsilon \in L(G)$, potom sa štandardne do množiny neterminálov pridáva nový začiatočný neterminál \hat{S} a do pravidiel \hat{P} sa pridajú 2 pravidlá:
 - $\hat{S} \rightarrow S$
 - $\hat{S} \rightarrow \varepsilon$
- To zistíme o.i. tak, že $S \in N_\varepsilon$.



Odstránenie ε -pravidiel - príklad

Príklad: Odstráňte ε -pravidlá z gramatiky $G = (\{S, A, B, C\}, \{a, b\}, P, S)$, kde pravidlá P :

$$S \rightarrow aAbB \mid AC$$

$$A \rightarrow Aa \mid \varepsilon$$

$$B \rightarrow bBb \mid ab$$

$$C \rightarrow aCA \mid \varepsilon \mid A$$

Množina N_ε postupne: $\emptyset, \{A, C\}, \{A, C, S\}$.



Odstránenie ε -pravidiel - príklad (pokr.)

Výsledná gramatika bude $\hat{G} = (\{\hat{S}, S, A, B, C\}, \{a, b\}, \hat{P}, \hat{S})$, kde \hat{P} :

$$\hat{S} \rightarrow S \mid \varepsilon$$

$$S \rightarrow aAbB \mid AC \mid abB \mid A \mid C$$

$$A \rightarrow Aa \mid a$$

$$B \rightarrow bBb \mid ab$$

$$C \rightarrow aCA \mid A \mid aC \mid aA \mid a$$

pričom boli odstránené pravidlá $A \rightarrow \varepsilon, C \rightarrow \varepsilon$



Odstránenie jednoduchých pravidiel

Definícia

*Nech $G = (N, T, P, S)$ je bezkontextová gramatika. Potom pravidlá tvaru $A \rightarrow B$, kde $A, B \in N$ sa nazývajú **jednoduché pravidlá**.*



Odstránenie jednoduchých pravidiel

Pre každý neterminál A vytvoríme množinu N_A neterminálov, ktoré sa z A dajú odvodiť, t.j.:

$$N_A = \{B \mid A \Rightarrow^* B, B \in N\} \quad (4)$$



Odstránenie jednoduchých pravidiel - príklad

Nech je daná gramatika (bola výsledkom úpravy odstránenia ε -pravidiel):

$$\acute{S} \rightarrow S \mid \varepsilon$$

$$S \rightarrow aAbB \mid AC \mid abB \mid A \mid C$$

$$A \rightarrow Aa \mid a$$

$$B \rightarrow bBb \mid ab$$

$$C \rightarrow aCA \mid A \mid aC \mid aA \mid a$$

Pre jednotlivé neterminály sú množiny N_A (postupne):

$$\acute{S} : \{\acute{S}\}, \{\acute{S}, S\}, \{\acute{S}, S, A, C\} = N_{\acute{S}}$$

$$S : \{S\}, \{S, A, C\} = N_S$$

$$A : \{A\} = N_A$$

$$B : \{B\} = N_B$$

$$C : \{C\} = \{C, A\} = N_C$$



Použitá literatúra

Dedera, L': Počítačové jazyky a ich spracovanie.

Linz, P.: An Introduction to Formal Languages and Automata.

Molnár, L': Gramatiky a jazyky.

