

Transformace bezkontextových gramatik

V tomto cvičení se budeme zabývat základními transformacemi bezkontextových gramatik - BKG (též CFG, context-free grammar). Popsané postupy a myšlenky jsou jednoduché a vcelku přímočaré, leč začátečník se může na chvíli v záplavě písmenek ztratit. Stane-li se vám to, nezoufejte. Zkuste na věc přijít selským rozumem a uvidíte, že na tom nic není :).

Nechť G je bezkontextová gramatika daná čtvericí (N, T, P, S) , kde N je množina neterminálních symbolů, T množina terminálních symbolů, P přepisovací pravidla a S startovací symbol. Připomeňme si, že gramatika je bezkontextová, tvoří-li levou stranu všech přepisovacích pravidel právě jeden neterminální symbol, tj. pravidla jsou ve tvaru $X \rightarrow w$, kde X je neterminální symbol a w řetězec složený z terminálních i neterminálních symbolů.

Gramatika generuje nějaký jazyk, tj. množinu řetězců terminálních symbolů. Tato množina může být konečná, nekonečná nebo i prázdná. Postupně se budeme zabývat těmito otázkami: Je generovaný jazyk prázdný? Jsou v gramatice takové symboly, které by v ní být nemusely? Dá se něco udělat s přepisovacími pravidly, aby gramatika vypadala "lépe"?

Ukázka algoritmu

Minulý rok vznikly v rámci DP demonstrace pro uváděné algoritmy, můžete se na ně podívat v [dočasném úložišti](#), na převodu do portálu se zatím pracuje.

Nedostupné symboly

Začneme-li vytvářet ze startovacího symbolu S gramatiky G větné formy, může se stát, že k některým symbolům (terminálním či neterminálním) vůbec nedojdeme. Takové symboly nazveme nedostupné. Základem algoritmu k jejich vyhledání je přímočarý postup, kterým "prolezáme" pravidla a sestpisujeme symboly, na které jsme narazili.

Můžete se podívat na [demonstraci na webu](#).

Algoritmus

- 1) $V\emptyset = \{S\}$, $i=1$
- 2) $V_i = \{X: A \rightarrow a \text{ a } b \text{ je přepisovací pravidlo, } a \text{ a } b \text{ jsou řetězce symbolů, } A \text{ náleží } V_{i-1}\} + V_{i-1}$
- 3) Je-li $V_i \neq V_{i-1}$, pak $i++$ a jdě na (2), jinak
 $VD = V_i$ (dostupné symboly)
 $VN = N + T - VD$ (nedostupné symboly)

Gramatiku bez nedostupných symbolů zapíšeme snadno - budeme psát jen ta přepisovací pravidla, ve kterých figurují dostupné symboly. Zároveň upravíme i množiny N a T.

Příklady

1. $G = (N, T, P, S)$, $N = \{S, A, B\}$, $T = \{a, b\}$, $P:$

$S \rightarrow aB \mid bB$
 $A \rightarrow a \mid BA$
 $B \rightarrow aS \mid bB$

Postupujeme podle algoritmu:

$V\emptyset = \{S\}$
 $V_1 = \{S, B, a, b\}$
 $V_2 = \{S, B, a, b\} = V_1$, tj. $VD = \{S, B, a, b\}$

Potom $G' = \{N', T', P', S\}$, kde $N' = \{S\}$, $T' = \{a, b\}$ a $P:$

$S \rightarrow aB \mid bB$
 $B \rightarrow aS \mid bB$

2. $G = (N, T, P, S)$, $N = \{S, A, B, C\}$, $T = \{a, b\}$, $P:$

$S \rightarrow A \mid aB$
 $A \rightarrow a \mid BA$
 $B \rightarrow bB \mid e$
 $C \rightarrow ABC \mid S$

Postupujeme podle algoritmu:

$V\emptyset = \{S\}$
 $V_1 = \{S, A, B, a\}$
 $V_2 = \{S, A, B, a, b\}$
 $V_3 = V_2$, tj. $VD = \{S, A, B, a, b\}$

Potom $G' = \{N', T', P', S\}$, kde $N' = \{S\}$, $T' = \{a, b\}$ a $P:$

$S \rightarrow A \mid aB$
 $A \rightarrow a \mid BA$
 $B \rightarrow bB \mid e$

Je jazyk prázdný?

Díky předchozímu textu umíme vyloučit z gramatiky takové symboly, které v ní nejsou evidentně užitečné. Za chvíli si ukážeme, že existuje ještě jeden typ takových symbolů. K tomu budeme potřebovat vědět, zda jazyk generovaný gramatikou G není prázdný, respektive budeme potřebovat algoritmus, který to zjistí.

Jazyk $L(G)$ je množina řetězců terminálních symbolů. Gramatika může mít ale takovou vlastnost, že při derivacích (přepisování) nebudou ubývat neterminální symboly. Nyní tedy chceme zjistit, zda z gramatiky G vůbec lze takový řetězec generovat.

Můžete se podívat na [demonstraci na webu](#).

Algoritmus:

- 1) $N\emptyset = \{\}$, $i = 1$

- 2) $N_i = \{A : A \rightarrow a, a \text{ je řetězec složený ze symbolů } N_{i-1} \text{ a } T\} + N_{i-1}$
- 3) Je-li $N_i \neq N_{i-1}$, i++ a jdi na (2), jinak
 $N_t = N_i$
- 4) Je-li S prvkem N_t , je $L(G)$ neprázdný, v opačném případě je prázdný.

Základní myšlenka je prostá. Pravidla bývají rekurzivní. Aby se mohl generovat řetězec terminálních symbolů, musí se rekurze někde zastavit. Nejprve najdeme neterminální symboly, u kterých se rekurze může zastavit díky vhodnému přepisovacímu pravidlu (symbol se přepisuje na prázdný řetězec nebo na řetězec terminálních symbolů). Takovým symbolům budeme říkat "dobré". Má-li nějaký jiný neterminální symbol X na pravé straně přepisovacího pravidla řetězec z terminálních symbolů a "dobrých" neterminálních symbolů, je i X "dobrý". A tak dále. Je-li "dobrý" i startovací symbol, je celá gramatika "dobrá", tj. generuje neprázdný jazyk. Jasné?

Příklady

3. gramatika G z příkladu 1 bez nedostupných symbolů, tj. $G = (N, T, P, S)$, $N = \{S, B\}$, $T = \{a, b\}$, P :

$$\begin{array}{l} S \rightarrow aB \mid bS \\ B \rightarrow aS \mid bB \end{array}$$

Postupujeme podle algoritmu:

$$\begin{array}{l} N_0 = \{\} \\ N_1 = N_0, \text{ protože } P \text{ neobsahuje žádné "dobré" pravidlo} \end{array}$$

Jazyk $L(G)$ je tedy prázdný.

4. $G = (N, T, P, S)$, $N = \{S, A, B\}$, $T = \{a, b\}$, P :

$$\begin{array}{l} S \rightarrow A \mid aB \\ A \rightarrow a \mid BA \\ B \rightarrow bB \mid e \end{array}$$

Postupujeme podle algoritmu:

$$\begin{array}{l} N_0 = \{\} \\ N_1 = \{A, B\} \\ N_2 = \{S, A, B\} = N_t \end{array}$$

Startovací symbol je "dobrý", tj. náleží do množiny N_t , tj. $L(G)$ je neprázdný.

5. $G = (N, T, P, S)$, $N = \{S, A, B, C, D\}$, $T = \{a, b\}$, P :

$$\begin{array}{l} S \rightarrow aAB \mid A \\ A \rightarrow bC \mid D \\ B \rightarrow A \mid b \\ C \rightarrow S \end{array}$$

Připadá-li vám, že je v zadání chyba, neboť pravidla ne definují přepisovací pravidla pro symbol D , dáváte pozor. Ale ne úplně, protože není nikde řečeno, že se všechny neterminální symboly musí vyskytovat na levé straně! Asi cítíte, že symbol D je v této gramatice zbytečný, tím se ale budeme zabývat až za chvíli. Nyní jen zjistíme, zda je $L(G)$ prázdný nebo ne.

$$\begin{array}{l} N_0 = \{\} \\ N_1 = \{B\} \\ N_2 = \{B\} = N_1 = N_t \end{array}$$

Nejenže gramatika obsahuje zbytečný symbol, gramatika je nějak zbytečná celá, protože generuje prázdný jazyk.

Zbytečné symboly

Víme, že některé symboly gramatiky mohou být nedostupné. Ty jsou zbytečné automaticky. V posledním příkladu jsme viděli, že symbol D byl také zbytečný, protože pro něj neexistovalo přepisovací pravidlo. Pokud si celý problém formalizujeme, zjistíme, že přepisování probíhá asi takto:

$$S \rightarrow^* a X b \rightarrow^* a c b$$

kde a , b a c jsou řetězce terminálních symbolů, X je neterminální symbol a S startovací symbol. Šipka s hvězdičkou představuje přepsání, které může být sérií základních derivací. Pokud pro symbol X neexistuje takové přepsání, je symbol X zbytečný. Jinými slovy, jakmile se ve větné formě objeví X , nikdy už z ní nevyloučí neterminální symboly.

Odstranění zbytečných symbolů probíhá ve dvou krocích. Nejprve najdeme množinu N_t . Nová gramatika bude obsahovat jen tyto "dobré" neterminální symboly. Z nich je ale zapotřebí vyházet všechny nedostupné symboly. Snadné, že?

Příklady

6. $G = (N, T, P, S)$, $N = \{S, A, B\}$, $T = \{a, b\}$, P :

$$\begin{array}{l} S \rightarrow ab \mid AB \\ A \rightarrow AB \\ B \rightarrow b \end{array}$$

Nejprve najdeme N_t :

$$\begin{array}{l} N_0 = \{\} \\ N_1 = \{S, B\} \\ N_2 = N_1 = N_t \end{array}$$

Pomocná gramatika je:

$$\begin{array}{l} S \rightarrow ab \\ B \rightarrow b \end{array}$$

Vyloučíme nedostupné symboly:

$$\begin{array}{l} V_0 = \{S\} \\ V_1 = \{S, a, b\} = VD \end{array}$$

Gramatika bez zbytečných symbolů obsahuje jen jedno pravidlo, a to:

$S \rightarrow ab$

7. $G = (N, T, P, S)$, $N = \{S, A, B, C, D\}$, $T = \{a, b, c\}$, $P:$

$S \rightarrow aB \mid bD$
 $A \rightarrow aB \mid bCD$
 $B \rightarrow ABCD$
 $C \rightarrow cC \mid A$
 $D \rightarrow aD \mid e$

Najdeme N_t :

$N_0 = \{\}$
 $N_1 = \{D\}$
 $N_2 = \{S, D\}$
 $N_3 = N_2 = N_t$

Nová gramatika má dvě pravidla (všechny symboly jsou dostupné):

$S \rightarrow bD$
 $D \rightarrow aD \mid e$

Vyloučení pravidla

V tomto okamžiku umíme z gramatiky vyloučit všechny nepotřebné symboly. Někdy je také vhodné upravit či vyloučit nějaká pravidla. Princip je prostý: místo toho, aby se přepisování provedlo při derivaci větné formy, provedeme ho už v zápisu gramatiky. Potřebujeme-li pozměnit pravidlo $A \rightarrow aB$ a všechna pravidla pro B jsou $B \rightarrow c_1, c_2, \dots, c_n$, přepíšeme pravidlo pro A na $A \rightarrow ac_1b, ac_2b, \dots, ac_nb$ (a, b, c jsou řetězce terminálních symbolů)

Speciálními pravidly, kterých se někdy potřebujeme zbavit, jsou tzv. e-pravidla, tj. přepsání neterminálního symbolu na prázdný řetězec. Gramatikou bez e-pravidel pak budeme rozumět takovou gramatiku, ve které nejsou žádná e-pravidla, až na jednu výjimku. Tou je pravidlo $s \rightarrow e$, které do gramatiky zahrneme v případě, že gramatika generuje prázdný řetězec (S je startovací symbol). Potom se ovšem nesmí symbol S objevit na žádné z pravých stran přepisovacích pravidel.

Algoritmus:

- 1) Způsobem podobným jako při hledání N_t najdeme N_e , množinu symbolů, které se dají přepsat na e .
- 2) Pravidlo $A \rightarrow a_0 B_1 a_1 B_2 a_2 \dots B_n a_n$ nahradíme pravidly $A \rightarrow a_0 X_1 a_1 X_2 a_2 \dots X_n a_n$, kde X_i je buď e nebo B_i (vyštědíme všech 2^i možností). Pravidlo typu $A \rightarrow e$ pochopitelně zapisovat nebudeme, protože takových pravidel se snažíme zbavit.
- 3) Patří-li startovací symbol S do množiny N_e , zařadíme do gramatiky pravidlo $S' \rightarrow S \mid e$. Novým startovacím symbolem pak bude S' .

Příklady

8. $G = (N, T, P, S)$, $N = \{S, A\}$, $T = \{a, b\}$, $P:$

$S \rightarrow aSbS \mid aSbA \mid A$
 $A \rightarrow aA \mid e$

Najdeme množinu N_e :

$N_0 = \{\}$
 $N_1 = \{A\}$ (protože A se dá přepsat na e)
 $N_2 = \{A, S\}$ (protože S se dá přepsat na A)
 $N_e = N_2$

Nyní budeme postupně vynechávat v pravidlech ty symboly, které se dají přepsat na e :

$S \rightarrow aSbS \mid aSb \mid abS \mid ab \mid$
 $\quad aSbA \mid abA \mid$
 $\quad A$
 $A \rightarrow aA \mid a$
 $S' \rightarrow S \mid e$ (protože S patří do N_e)

9. $G = (N, T, P, S)$, $N = \{S, A, B\}$, $T = \{a, b\}$, $P:$

$S \rightarrow aB \mid bA$
 $A \rightarrow aA \mid B \mid e$
 $B \rightarrow bB \mid A \mid e$

Najdeme množinu N_e :

$N_0 = \{\}$
 $N_1 = \{A, B\}$
 $N_2 = N_1 = N_e$

Po přepisu:

$S \rightarrow aB \mid bA \mid a \mid b$
 $A \rightarrow aA \mid a \mid B$
 $B \rightarrow bB \mid b \mid A$

Jednoduchá pravidla

Jednoduché pravidlo ja takové, při kterém se neterminální symbol přepisuje na jiný neterminální symbol, tedy pravidlo typu $A \rightarrow B$. Sama o sobě by taková pravidla v gramatice nevadila, i když se dá diskutovat nad tím, zda jsou skutečně nutná. Hlavní motivací jejich odstranění je vedlejší jev, a to cykly. Na gramatice v posledním příkladu můžeme takový cyklus vidět: A se přepisuje na B a to opět na A . Odstraněním jednoduchých pravidel tedy odstraníme i cykly, které jsou nežádoucí při syntaktické analýze.

Nejprve pro každý neterminální symbol zjistíme, na jaké neterminální symboly se může přepsat pomocí jednoduchých pravidel. To se dá

zjistit postupem podobným jako při hledání dostupných symbolů. Pro neterminální symbol A zařadíme do množiny N_A jak symbol A, tak všechny "jednoduše dostupné" symboly.

V dalším kroku již píšeme nová pravidla. Pro neterminální symbol A napíšeme vechna pravidla $A \rightarrow a$, kde a jsou ne-jednoduchá pravidla pro všechny symboly z N_A .

Příklady

10. $G = (N, T, P, E)$, $N = \{E, T, F\}$, $T = \{+, *, (,), i\}$, P :

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid i \end{array}$$

V této gramatice se sice cykly nevyskytují, ale z cvičných důvodů jednoduchá pravidla odstraníme:

$$\begin{array}{l} NE = \{E, T, F\} \\ NT = \{T, F\} \\ NF = \{F\} \end{array}$$

Nová gramatika je tedy:

$$\begin{array}{l} E \rightarrow E + T \mid T * F \mid (E) \mid i \\ T \rightarrow T * F \mid (E) \mid i \\ F \rightarrow (E) \mid i \end{array}$$

11. $G = (N, T, P, S)$, $N = \{A, B, S\}$, $T = \{a, b\}$, P :

$$\begin{array}{l} S \rightarrow aB \mid bA \\ A \rightarrow aA \mid a \mid B \\ B \rightarrow bB \mid b \mid A \end{array}$$

Sestavíme množiny $N_?$:

$$\begin{array}{ll} NS = \{S\} & (\text{zajímají nás jen jednoduchá pravidla!}) \\ NA = \{A, B\} \\ NB = \{A, B\} \end{array}$$

Nová gramatika je:

$$\begin{array}{l} S \rightarrow aB \mid bA \\ A \rightarrow a \mid aA \mid b \mid bB \\ B \rightarrow a \mid aA \mid b \mid bB \end{array}$$

Vlastní gramatika je gramatika bez zbytečných symbolů, bez e-pravidel a bez cyklů.

Odstranění levé rekurze

Pravidlo je levorekurzivní, je-li ve tvaru $A \rightarrow A \dots$. Gramatiky s levorekurzivními pravidly jsou nevhodné pro některé typy syntaktické analýzy. Jsou-li všechna pravidla pro neterminální symbol A zapsána jako

$$A \rightarrow Aa_1 \mid Aa_2 \mid \dots \mid Aa_m \mid b_1 \mid b_2 \mid \dots \mid b_n$$

kde a_i a b_i jsou řetězce terminálních a neterminálních symbolů a žádné b_i nezačíná symbolem A, můžeme odstranit přímou levou rekurzi takto:

$$\begin{array}{l} A \rightarrow b_1 \mid b_2 \mid \dots \mid b_n \mid b_1A' \mid b_2A' \mid \dots \mid b_nA' \\ A' \rightarrow a_1 \mid a_2 \mid \dots \mid a_m \mid a_1A' \mid a_2A' \mid \dots \mid a_mA' \end{array}$$

nebo pomocí e-pravidel:

$$\begin{array}{l} A \rightarrow b_1A' \mid b_2A' \mid \dots \mid b_nA' \\ A' \rightarrow e \mid a_1A' \mid a_2A' \mid \dots \mid a_mA' \end{array}$$

Sice jsme se učili e-pravidla odstraňovat, ale někdy jsou menší zlo.

Nepřímá levá rekurze je například toto: $A \rightarrow B \dots$, $B \rightarrow A \dots$. Její odstranění funguje podobně jako u přímé rekurze. Je velmi vhodné si do odstraňování zavést nějaký systém, jinak je slušná šance, že se v záplavě písmenek člověk ztratí!

Jeden ze systémů je následující. Označíme si neterminální symboly jako N_1, N_2, \dots, N_n . Jsou-li v gramatice pravidla typu $N_i \rightarrow N_j \dots$, budeme se snažit, aby pokud možno $i < j$, tj. aby se přepisování neodkazovalo na předcházející symboly. Není to nutné, ale může to ušetřit práci. Pro pravidlo $N_i \rightarrow \dots$ nejprve zajistíme, aby se neodkazovalo na předcházející symboly. I v tom je třeba zavést pořádek, proto nejdříve odstraníme odkaz na N_1 , pak na N_2 atd. Až na závěr si necháme odstranění přímé levé rekurze $N_i \rightarrow N_i$. Toto "čištění" pravidel provedeme postupně pro N_1, N_2 atd.

Příklady

12. $G = (N, T, P, E)$, $N = \{E, T, F\}$, $T = \{+, *, (,), i\}$, P :

$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid i \end{array}$$

Vidíme, že zde nepřímá levá rekurze není, budeme mít méně práce. Nejprve bez e-pravidel:

$$\begin{array}{l} E \rightarrow T \mid TE' \\ E' \rightarrow + T \mid + TE' \\ T \rightarrow F \mid FT' \\ T' \rightarrow * F \mid * FT' \\ F \rightarrow (E) \mid i \end{array}$$

V dalším cvičení zjistíme, že pro syntaktickou analýzu je vhodnější zavedení e-pravidel:

```
E --> TE'  
E' --> e | + TE'  
T --> FT'  
T' --> e | * FT'  
F --> (E) | i
```

13. $G = (N, T, P, A)$, $N = \{A, B, C\}$, $T = \{a, b\}$, P :

```
A --> AB | C | Bb  
B --> Ab | Ca  
C --> Aa | b
```

Zde se nepřímá levá rekurze vyskytuje. Při uspořádání symbolů A-B-C se B i C odkazují na předešlý symbol, A. V uspořádání B-C-A se na předešlý symbol odkazuje jen A. Použijeme tedy pořadí B-C-A.

```
B --> Ab | Ca      (první pravidlo se nemůže odkazovat na předešlé symboly,  
                      přímá levá rekurze zde není)  
C --> Aa | b       (neodkazuje se na B, přímá levá rekurze zde není)  
  
A --> AB | C | Bb          (původní pravidla)  
A --> AB | C | Abb | Cab    (odstranili jsme odkaz na B)  
A --> AB | Aa | b | Abb | Aaab | bab   (odstranili jsme odkaz na C)  
  
A --> b | bab | bA' | babA'     (nyní už píšeme definitivní pravidla)  
A' --> B | a | bb | aab | BA' | aA' | bbA' | aabA'
```

S e-pravidly:

```
B --> Ab | Ca  
C --> Aa | b  
  
A --> AB | C | Bb          (původní pravidla)  
A --> AB | C | Abb | Cab    (odstranili jsme odkaz na B)  
A --> AB | Aa | b | Abb | Aaab | bab   (odstranili jsme odkaz na C)  
  
A --> bA' | babA'          (definitivní pravidla)  
A' --> e | BA' | aA' | bbA' | aabA'
```

Příklady na domácí procvičení (vyřešte nejméně 3)

1. Převeďte $G = (N, T, P, S)$ na gramatiku bez zbytečných symbolů:

```
S --> A | B  
A --> aB | bS | b  
B --> AB | Ba  
C --> AS | b
```

2. Odstraňte u $G=(N, T, P, S)$ e-pravidla:

```
S --> ABC  
A --> BB | e  
B --> CC | a  
C --> AA | b
```

3. Odstraňte u $G=(N, T, P, S')$ jednoduchá pravidla:

```
S' --> S | e  
S --> A | B  
A --> C | D  
B --> D | E  
C --> S | a  
D --> S | b  
E --> S | c
```

4. Odstraňte u $G=(N, T, P, A)$ levou rekurzi bez použití e-pravidel:

```
A --> Ba | c  
B --> CA | bb  
C --> Ab | Bb
```

5. Odstraňte u $G=(N, T, P, S)$ levou rekurzi s použitím e-pravidel:

```
S --> AB  
A --> BS | b  
B --> SA | a
```

