

Gramatyki i języki bezkontekstowe \mathcal{L}_2

Produkcje bezkontekstowe: $X \rightarrow \omega, \quad X \in V, \quad \omega \in (A \cup V)^*$

Produkcje kontekstowe: $\alpha X \beta \rightarrow \alpha \omega \beta, \quad \omega \neq \lambda$

Przykłady

• $L = \{a^n b^n : n \geq 0\} \quad S \rightarrow aSb \mid \lambda$

• $L = \{a^m b^n : m \neq n\}$

$L = L_1 \cup L_2$, gdzie $L_1 = \{a^m b^n : m > n\}$, $L_2 = \{a^m b^n : m < n\}$.

$$S_0 \rightarrow AS \mid SB$$

$$S \rightarrow aSb \mid \lambda$$

$$A \rightarrow aA \mid a$$

$$B \rightarrow bB \mid b.$$

• $L = \{\omega \in A^* : |\omega|_a = |\omega|_b\} \quad S \rightarrow SS \mid aSb \mid bSa \mid \lambda$

$$S \rightarrow SS \mid (S) \mid \lambda$$

• Język prostych wyrażeń arytmetycznych

$$S \rightarrow S + S \mid S * S \mid X$$

$$X \rightarrow x \mid y \mid z \mid (S)$$

Przeanalizować wyprowadzenia

$$aaabbbbbb, \quad aaaaabbb, \quad (()()), \quad (()())(), \quad x * (y + z), \quad x * y + z$$

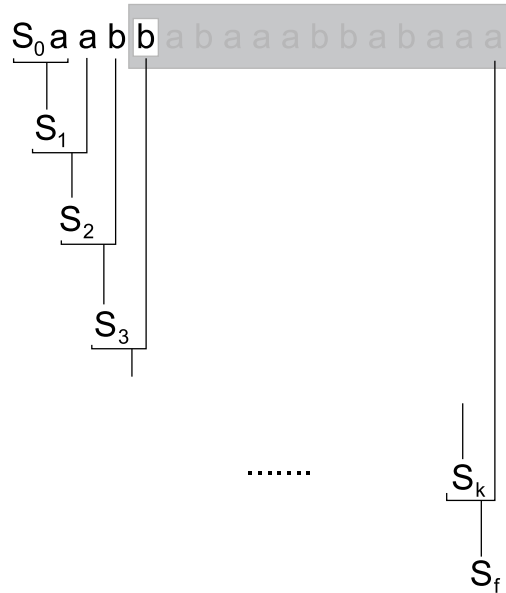
Przykłady:

$$L = \{a^n b^m : 2n \leq m \leq 3n\}, \quad L = \{\omega : |\omega|_a \neq |\omega|_b\}, \quad L = \{\omega \overleftarrow{\omega} : \omega \in A^*\}^c$$

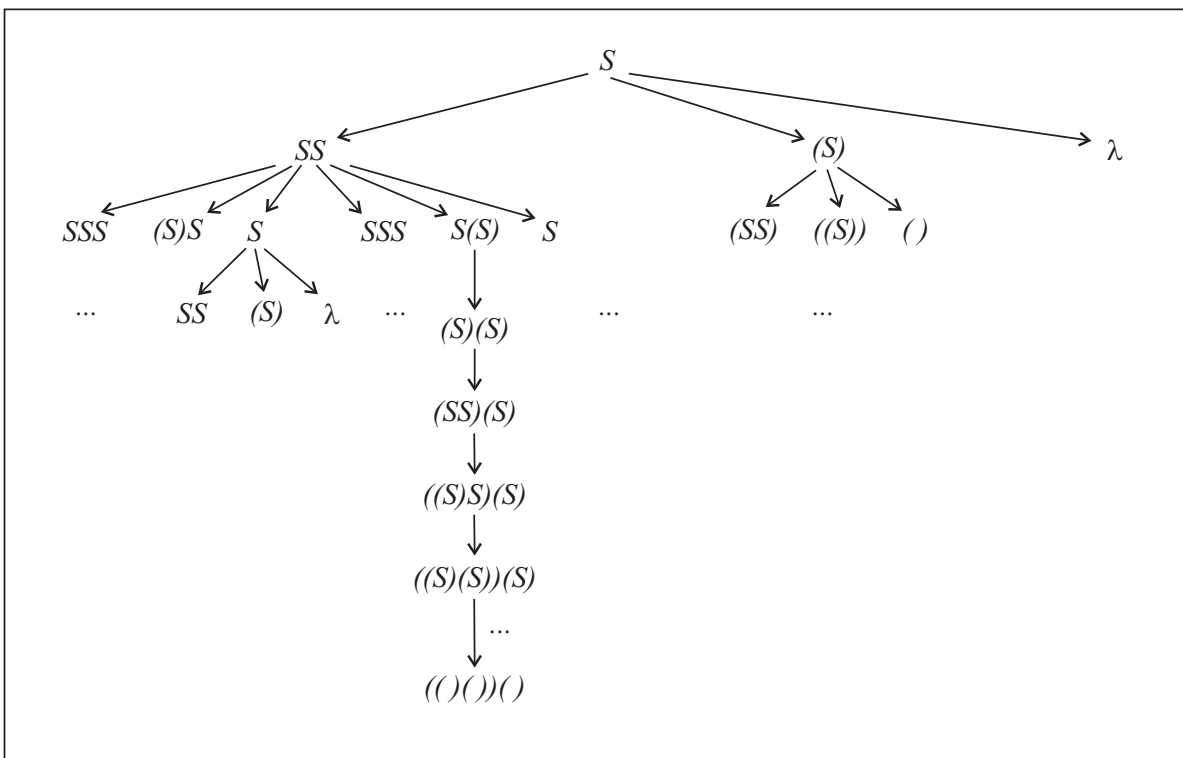
Jak sprawdzić mając daną gramatykę czy dane słowo $\alpha \in L(G)$?

a a a b b b b b

Analiza bottom-up i top-down



$$S \rightarrow SS \mid (S) \mid \lambda$$



Szacowanie złożoności przeszukania drzewa wyprowadzeń

- gramatyka bez pustych i jednostkowych produkcji: $X \rightarrow \lambda$, $X \rightarrow Y$
- $X \rightarrow \beta$, $|\beta| > 1$ za wyjątkiem $X \rightarrow a$ (produkcje nieskracające)
- liczba wszystkich produkcji gramatyki p

Wówczas

- liczba węzłów na k -tym piętrze drzewa wyprowadzeń $\leq p^k$
- wszystkie $\alpha \in A^*$ długości $\leq n$ znajdują się w piętrach $\leq 2n$
- przeszukanie drzewa by sprawdzić czy $\alpha \in L(G)$ (o długości $\leq n$)

$$1 + p + p^2 + \dots + p^{2n} = O(p^{2n+1})$$

By usprawnić wyprowadzanie napisów, szukamy dogodnej postaci gramatyki:

- bez pustych i jednostkowych produkcji
- bez wielu różnych wyprowadzeń tego samego napisu
- uporządkowanie procesu wyprowadzania: wyprowadzenia lewostronne

Jeszcze raz język prostych wyrażeń arytmetycznych

$$\begin{array}{ll} S \rightarrow S + S \mid S * S \mid X & S \rightarrow T \mid T + S \\ X \rightarrow x \mid y \mid z \mid (S) & T \rightarrow F \mid F * T \\ & F \rightarrow x \mid y \mid z \mid (S) \end{array}$$

$$S \Rightarrow S + S \Rightarrow S + S * S \Rightarrow S + S * S + S \Rightarrow \dots$$

$$S \Rightarrow S * S \Rightarrow S + S * S \Rightarrow S + S * S + S \Rightarrow \dots$$

$$S \Rightarrow S + S \Rightarrow S + S + S \Rightarrow S + S * S + S \Rightarrow \dots$$

W drugiej gramatyce

$$\begin{aligned} S &\Rightarrow T + S \Rightarrow F + S \Rightarrow x + S \Rightarrow x + T + S \Rightarrow x + F * T + S \\ &\Rightarrow x + y * T + S \Rightarrow x + y * F + S \Rightarrow x + y * z + S \Rightarrow x + y * z + T \\ &\Rightarrow x + y * z + F \Rightarrow x + y * z + z. \end{aligned}$$

Definicja Mówimy że gramatyka jest jednoznaczna, jeśli dla dowolnego napisu $\alpha \in A^*$ istnieje co najwyżej jedno wyprowadzenie lewostronne z S .

Istnieją języki bezkontekstowe wewnątrznie niejednoznaczne

Postać normalna gramatyki bezkontekstowej

Definicja Mówimy, że bezkontekstowa gramatyka G jest dana w normalnej postaci Chomsky'ego, jeśli wszystkie jej produkcje mają formę

$$X \rightarrow YZ \quad \text{lub} \quad X \rightarrow a,$$

gdzie $X, Y, Z \in V$ oraz $a \in A$.

Twierdzenie Każdą gramatykę bezkontekstową $G = (A, V, S, \Pi)$ taką, że $\lambda \notin L(G)$, można przekształcić do równoważnej postaci normalnej Chomsky'ego $G' = (A, V', S, \Pi')$.

Przekształcenie gramatyki do postaci normalnej Chomsky'ego

- 1) Usuwamy z gramatyki produkcje puste.
- 2) Eliminujemy produkcje jednostkowe.
- 3) Usuwamy symbole i produkcje bezużyteczne.
- 4) Jeśli $\lambda \in L(G)$, uzupełniamy G' o dodatkową produkcję $S \rightarrow \lambda$. Jest to jedyna pusta produkcja dopuszczalna w G' .
- 5) Przekształcenie produkcji gramatyki do postaci normalnej Chomsky'ego:

$$X \rightarrow YZ, \quad X \rightarrow a.$$

Nagrania mp4 z przykładami redukcji gramatyk do postaci Chomsky'ego – p. strona przedmiotu, link z datą 18.XII.20.

Algorytm Cooke'a, Youngera i Kasamiego (CYK): dla zadanej gramatyki bezkontekstowej G w postaci normalnej Chomsky'ego i słowa α zbadać czy $\alpha \in L(G)$.

Niech $\alpha = a_1 a_2 \dots a_n$. Oznaczmy przez $\alpha_{ij} = a_i a_{i+1} \dots a_j$. Działanie algorytmu polega na stopniowym tworzeniu zbiorów symboli nieterminalnych

$$V_{ij} = \{X \in V : X \Rightarrow^* \alpha_{ij}\},$$

kolejno od

$$V_{ii} = \{X \in V : X \rightarrow a_i \text{ jest produkcją w } G\}$$

aż do zbioru

$$V_{1n} = \{X \in V : X \Rightarrow^* \alpha_{1n} = \alpha\}.$$

Wystarczy wówczas sprawdzić czy $S \in V_{1n}$.

Niech $X \in V_{ij}$, $i < j$, czyli $X \Rightarrow^* \alpha_{ij} = a_i a_{i+1} \dots a_j$.

Wyprowadzenie to zaczyna się od zastosowania pewnej produkcji $X \rightarrow YZ$,

$$X \Rightarrow YZ \Rightarrow^* \alpha_{ij} = a_i a_{i+1} \dots a_k a_{k+1} \dots a_j,$$

czyli

$$Y \Rightarrow^* a_i a_{i+1} \dots a_k = \alpha_{ik} \quad \text{oraz} \quad Z \Rightarrow^* a_{k+1} \dots a_j = \alpha_{k+1j}$$

dla pewnego $i \leq k < j$. Stąd $Y \in V_{ik}$ i $Z \in V_{k+1j}$. Odwrotnie

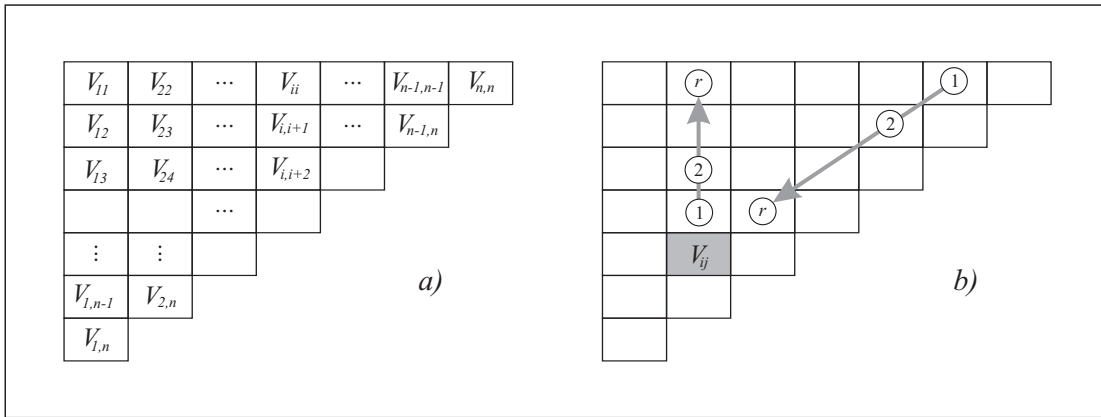
$$X \in V_{ij} \quad \Leftrightarrow \quad \exists i \leq k < j \quad X \rightarrow YZ \in \Pi, \quad Y \in V_{ik} \quad \text{i} \quad Z \in V_{k+1j}.$$

Konstrukcja zbiorów V_{ij}

$$\begin{array}{ccc} Y \in V_{ik} & & V_{k+1j} \ni Z \\ & \searrow & \swarrow \\ & V_{ij} & \\ & X \rightarrow YZ & \end{array}$$

kolejno dla $k = i, i+1, \dots, j-1$. A więc

$$V_{ij} = \bigcup_{i \leq k < j} \{X : X \rightarrow YZ \in \Pi \quad \text{oraz} \quad Y \in V_{ik}, \quad Z \in V_{k+1j}\}$$



Przykład. Weźmy jako przykład gramatykę $S \rightarrow () | (S) | SS$ po przekształceniu jej do postaci Chomsky'ego

$$\begin{aligned}
 S &\rightarrow LR | LQ | SS \\
 Q &\rightarrow SR \\
 L &\rightarrow (\\
 R &\rightarrow)
 \end{aligned}$$

Tworzenie pierwszego wiersza tabelki V_{ij} dla napisu $((()()))$

(()	())	()
L	L	R	L	R	R	L	R

Tworzenie drugiego wiersza: $V_{i,i+1}$ to zbiór tych symboli X , które są lewymi stronami produkcji $X \rightarrow YZ$, gdzie $X \in V_{ii}$, a $Y \in V_{i+1,i+1}$.

(()	())	()
L	L	R	L	R	R	L	R
\emptyset	S	\emptyset	S	\emptyset	\emptyset	S	

Pomocnicza tabela podglądu: $TP[Y, Z] = \{X \in V : X \rightarrow YZ \in \Pi\}$. Dla naszego przykładu

	S	L	R	Q
S	S	\emptyset	Q	\emptyset
L	\emptyset	\emptyset	S	S
R	\emptyset	\emptyset	\emptyset	\emptyset
Q	\emptyset	\emptyset	\emptyset	\emptyset

Końcowa postać tabeli

(() ()) ()

L	L	R	L	R	R	L	R
\emptyset	S	\emptyset	S	\emptyset	\emptyset	S	
\emptyset	\emptyset	\emptyset	Q	\emptyset	\emptyset		
\emptyset	S	\emptyset	\emptyset	\emptyset			
\emptyset	Q	\emptyset	\emptyset				
S	\emptyset	\emptyset					
\emptyset	\emptyset						
S							

Niech $\alpha = a_1 a_2 \dots a_n$ oraz niech tablica $TP[Y, Z]$ zawiera informację o lewych stronach produkcji $X \rightarrow YZ$. Prócz tego gramatyka zawiera produkcje terminalne $X_i \rightarrow a_i$. Pseudokod CYK:

```
function LP(U, W){
    V = empty;
    for (Y ∈ U)
        for (Z ∈ W)
            V = V ∪ TP[Y,Z];
    return V;
}

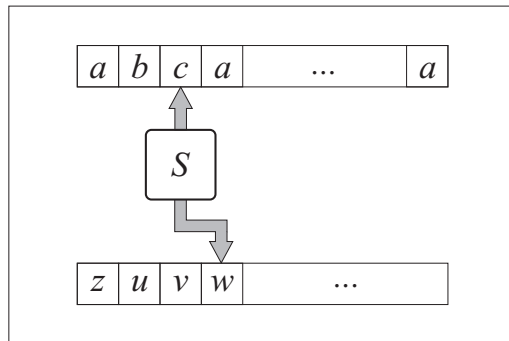
for (i=1, ..., n) Vii = {Xi};

for (l=1, ..., n-1)
    for (i=1, ..., n-1)
        j=i+1;
        Vij = empty;
        for (k=i, ..., j-1)
            Vij = Vij ∪ LP(Vik, Vk+1j);

if (S ∈ V1n) print("OK");
```

Złożoność wynosi $O(n^3)$ – potrójna pętla.

Automaty ze stosem



Definicja: *Niedeterministyczny automat ze stosem (NAS) opisany jest przez siedem obiektów*

$$M = (A, V, z, \Sigma, s_0, S_F, \pi),$$

gdzie:

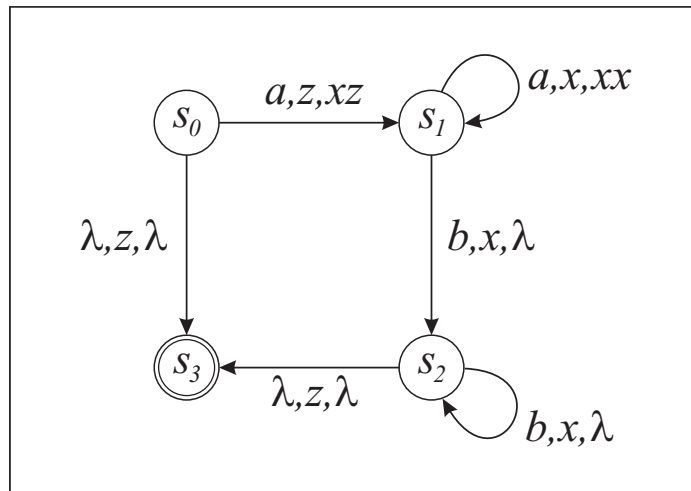
- A jest alfabetem wejściowym*
- V jest alfabetem stosu*
- $z \in V$ jest wyróżnionym symbolem początkowym stosu*
- Σ jest zbiorem stanów automatu, $\sigma = \{s_0, s_1, \dots, s_q\}$*
- $s_0 \in \Sigma$ jest wyróżnionym stanem startowym*
- $S_F \subseteq \Sigma$ jest zbiorem stanów końcowych (akceptujących)*
- π jest relacją przejścia, $\pi : \Sigma \times (A \cup \{\lambda\}) \times V \rightarrow \mathcal{F}(\Sigma \times V^*)$.*

Relacja przejścia $\pi(\text{stan}, \text{litera}, \text{symbol stosu}) = \text{zbiór par}(\text{stan}, \text{napis do stosu})$

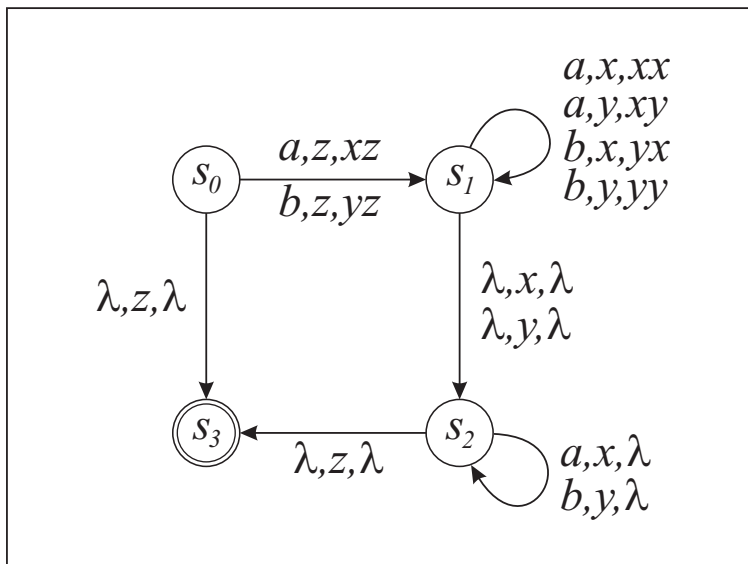
$$\pi(s_i, a, x) = \{(s_{j_1}, \gamma_1), (s_{j_2}, \gamma_2), \dots, (s_{j_k}, \gamma_k)\}.$$

- Jeśli czytana litera to a , stan automatu s_i oraz znak na wierzchu stosu to x , automat przesuwa się w prawo do następnej litery, zmienia stan na jeden z możliwych stanów s_j , usuwa znak x ze stosu, a w jego miejsce zapisuje odpowiadający s_j ciąg symboli stosu γ_j , być może λ .
- Przejście puste $\pi(s_i, \lambda, x) = (s_j, \gamma_j)$ – bez odczytywania kolejnej litery i ruchu głowicy czytającej.
- Jeśli stos jest pusty (nie ma x do odczytania), automat zatrzymuje się nie mogąc wykonać kolejnego ruchu (stąd specjalny symbol z jako początkowa zawartość stosu).
- Napis wejściowy zostanie zaakceptowany, jeśli po przeczytaniu ostatniej jego litery automat znajdzie się w jednym ze stanów końcowych. Zawartość stosu w tym momencie jest nieistotna.

Przykład: Automat akceptujący język $L = \{a^n b^n : n \geq 0\}$



Przykład: Automat akceptujący język $L = \{\alpha \bar{\alpha} : \alpha \in A^*\}$



Definicja: Automat ze stosem jest deterministyczny (DAS) jeśli relacja przejścia π jest częściową funkcją, a więc posiada następujące własności:

- (i) dla każdej trójki $(s, a, x) \in \Sigma \times A \times V$ istnieje co najwyżej jedna para $(s', \gamma) \in \Sigma \times V^*$ taka, że $\pi(s, a, x) = (s', \gamma)$;
- (ii) jeśli puste przejście $\pi(s, \lambda, x)$ jest określone dla pewnych s i x , wtedy $\pi(s, a, x)$ nie może być określone dla jakiegokolwiek litery a .

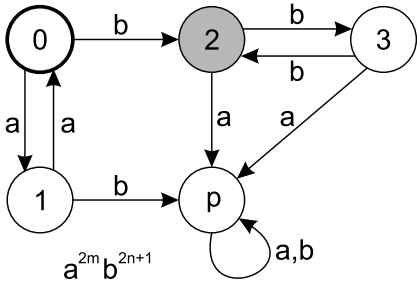
Twierdzenie: Klasa \mathcal{L}_2 języków bezkontekstowych (tj. opisywanych przez gramatyki bezkontekstowe) jest identyczna z klasą języków akceptowanych przez niedeterministyczne automaty ze stosem \mathcal{L}_{NAS} . Języki klasy \mathcal{L}_{DAS} stanowią właściwy podzbiór \mathcal{L}_2

$$\mathcal{L}_{\text{DAS}} \subsetneq \mathcal{L}_{\text{NAS}} = \mathcal{L}_2.$$

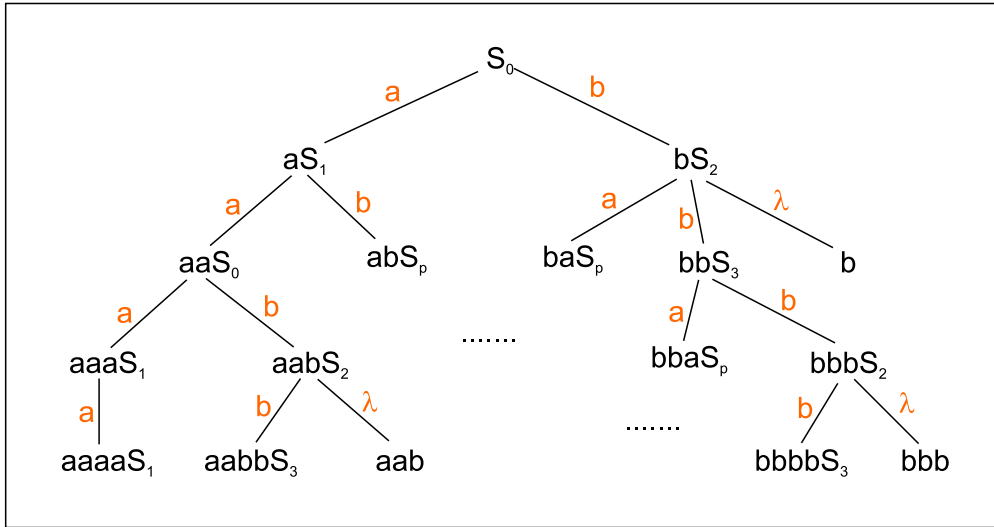
Przykład: Język niedeterministyczny

$$L = L_1 \cup L_2 = \{a^n b^n : n \geq 0\} \cup \{a^n b^{2n} : n \geq 0\}.$$

Po przeczytaniu ostatniej litery a nadal nie wiadomo czy przetwarzane słowo może pochodzić z L_1 czy z L_2 . Decyzja automatu co do wyboru dalszej ścieżki musi być podjęta niedeterministycznie.



- $S_0 \rightarrow aS_1 \mid bS_2$
- $S_1 \rightarrow aS_0 \mid bS_p$
- $S_2 \rightarrow aS_p \mid bS_3 \mid \lambda$
- $S_3 \rightarrow aS_p \mid bS_2$
- $S_p \rightarrow aS_p \mid bS_p$



Języki formalne i automaty 4.I.2022

Gramatyki i języki bezkontekstowe \mathcal{L}_2

Twierdzenie: Klasa \mathcal{L}_2 języków bezkontekstowych (tj. opisywanych przez gramatyki bezkontekstowe) jest identyczna z klasą języków akceptowanych przez nie-deterministyczne automaty ze stosem \mathcal{L}_{NAS} . Języki klasy \mathcal{L}_{DAS} stanowią właściwy podzbiór \mathcal{L}_2

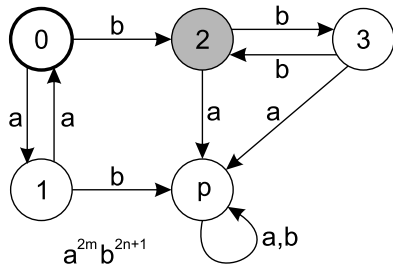
$$\mathcal{L}_{\text{DAS}} \subsetneq \mathcal{L}_{\text{NAS}} = \mathcal{L}_2.$$

Przykład: $L = L_1 \cup L_2 = \{a^n b^n : n \geq 0\} \cup \{a^n b^{2n} : n \geq 0\}$.

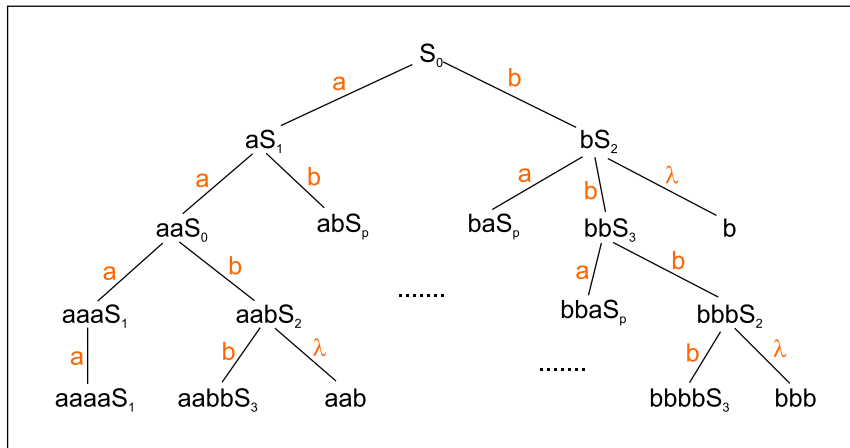
Po przeczytaniu ostatniej litery a nadal nie wiadomo czy przetwarzane słowo może pochodzić z L_1 czy z L_2 . Decyzja automatu co do wyboru dalszej ścieżki musi być podjęta niedeterministycznie.

Niedeterminizm jest źródłem niższej efektywności algorytmów weryfikacji czy $\alpha \in L$. W praktycznych zastosowaniach (składnia języków programowania) wybieramy struktury deterministyczne.

Przykład: Gramatyka regularna na bazie automatu deterministycznego



$$\begin{aligned} S_0 &\rightarrow aS_1 \mid bS_2 \\ S_1 &\rightarrow aS_0 \mid bS_p \\ S_2 &\rightarrow aS_p \mid bS_3 \mid \lambda \\ S_3 &\rightarrow aS_p \mid bS_2 \\ S_p &\rightarrow aS_p \mid bS_p \end{aligned}$$



Specjalna klasa gramatyk **s-gramatyki** (simple)

Definicja: G jest s-gramatyką, jeśli wszystkie jej produkcje mają postać

$$X \rightarrow a Y_1 \dots Y_m, \quad \text{gdzie } a \in A, \quad Y_i \in V, \quad m \geq 0,$$

przy czym dla każdej pary X i a istnieje co najwyżej jedna taka produkcja.

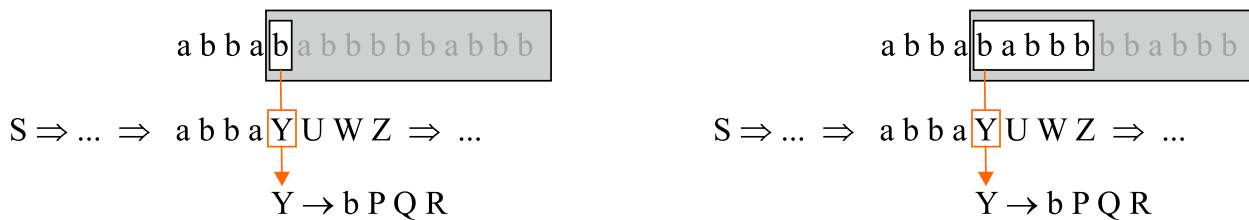
Produkcje s-gramatyki można zapisać w tablicy podglądu:

$$T[X, a] = Y_1 \dots Y_m \quad \text{lub} \quad T[X, a] = \emptyset$$

Niech $\alpha = a_1 a_2 \dots a_n$, tworzymy wyprowadzenie lewostronne zaczynając od symbolu S .

$$\begin{aligned} S &\Rightarrow a_1 Y_1 \dots Y_m \\ T[S, a_1] &= Y_1 \dots Y_m \\ &\Rightarrow a_1 a_2 Z_1 \dots Z_k Y_2 \dots Y_m \\ T[Y_1, a_2] &= Z_1 \dots Z_k \\ &\Rightarrow a_1 a_2 a_3 \dots \\ T[Z_1, a_3] &= \dots \end{aligned}$$

Podgląd jednego kolejnego symbolu wejściowego napisu pozwala na jednoznaczne ustalenie kolejnej produkcji w wyprowadzeniu (lub jej braku).



Definicja: Gramatyka bezkontekstowa $G = (A, V, S, \Pi)$ jest typu $LL(k)$, jeśli dla każdej pary wyprowadzeń lewostronnych postaci

$$\begin{aligned} S &\Rightarrow^* \alpha_1 X \omega_1 \Rightarrow \alpha_1 \beta_1 \omega_1 \Rightarrow^* \alpha_1 \alpha_2 \\ S &\Rightarrow^* \alpha_1 X \omega_2 \Rightarrow \alpha_1 \beta_2 \omega_2 \Rightarrow^* \alpha_1 \alpha_3, \end{aligned}$$

gdzie $\alpha_i \in A^*$, $\beta_i, \omega_i \in (A \cup V)^*$, równość co najwyżej k początkowych liter w α_2 i α_3 pociąga za sobą równość $\beta_1 = \beta_2$.

Przykłady

$G = \{S \rightarrow aSB|aB, B \rightarrow b\}$ nie jest s-gramatyką, bo obydwie produkcje dla S rozpoczynają się literą a . G jest gramatyką typu $LL(2)$:

	aa	ab	ba	bb
S	aSB	aB	\emptyset	\emptyset

Istnieje również s-gramatyka:

$$\begin{aligned} S &\rightarrow aA \\ A &\rightarrow aAB|b \\ B &\rightarrow b \end{aligned}$$

$G = \{S \rightarrow ()|(S)|SS\}$ nie jest typu $LL(k)$ dla żadnego k . Przy analizowaniu napisów o długości przekraczającej 2 mamy do wyboru dwie produkcje: $S \rightarrow (S)$ oraz $S \rightarrow SS$. Skanowany aktualnie symbol (ani k następnych) nie podpowiada, której z nich powinniśmy użyć. Weźmy bowiem słowa

$$((\dots() \dots)) \quad \text{albo} \quad ((\dots() \dots))().$$

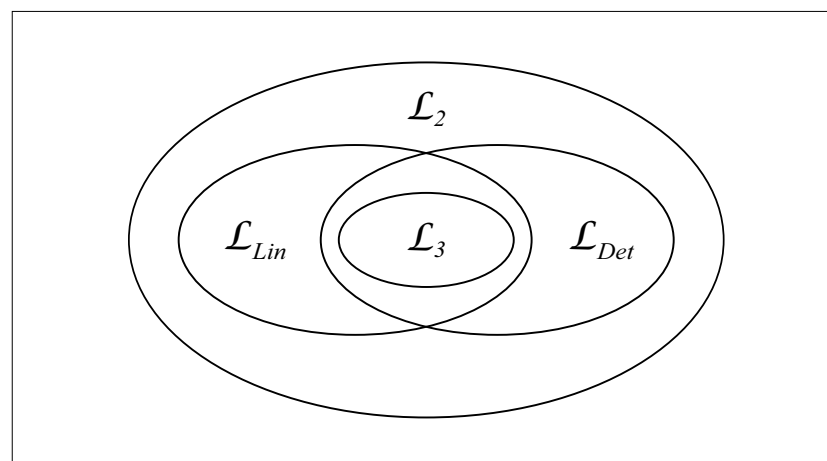
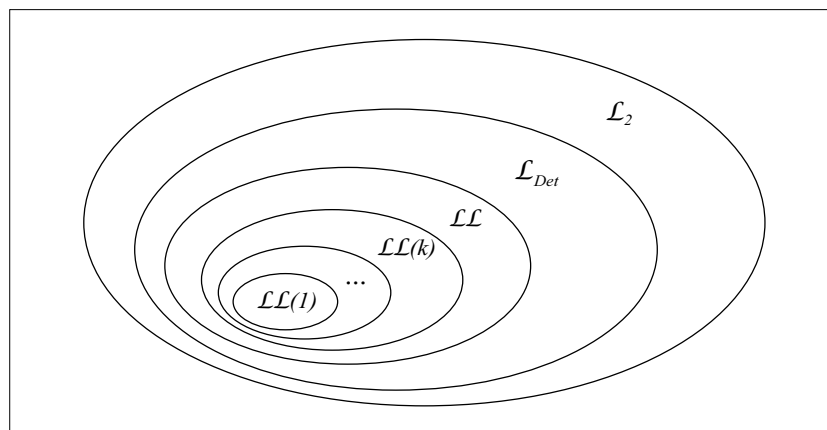
Wyprowadzenie pierwszego z nich polega na iteracyjnym użyciu produkcji $S \rightarrow (S)$, natomiast wyprowadzenie drugiego *musi* zaczynać się od $S \rightarrow SS$. Decyzję co do tego, której produkcji użyć w wyprowadzeniu jako pierwszej można podjąć na podstawie podglądu *nie mniej* niż $2n$ początkowych symboli wejściowych. Ponieważ jednak w naszej gramatyce możemy wyprowadzać takie napisy dla dowolnie dużych n , nie może być ona typu $LL(k)$ dla jakiegokolwiek ustalonego k .

Ten sam język posiada gramatykę typu $LL(1)$: $G = \{S \rightarrow (S)S|\lambda\}$. Wyprowadzenie $((\dots)())$:

$$S \xRightarrow{\lambda} (S)S \xRightarrow{\lambda} ((S)S)S \xRightarrow{\lambda} ((\dots)S)S \xRightarrow{\lambda} ((\dots)S) \xRightarrow{\lambda} ((\dots)S)S \xRightarrow{\lambda} ((\dots)())S \xRightarrow{\lambda} ((\dots)()).$$

Przeanalizować samodzielnie **Przykład 4.12** str. 105 z gramatyką $LL(1)$ dla wyrażeń arytmetycznych.

Klasy języków deterministycznych



Lematy o pompowaniu dla języków bezkontekstowych

Na początek języki liniowe:

$$X \rightarrow \alpha Y \beta, \quad X \rightarrow \alpha, \quad \alpha, \beta \in A^*$$

Zakładamy przy tym, że produkcje są nieskracające i niejednostkowe (α, β nie mogą być jednocześnie puste). Niech N będzie liczbą symboli nieterminalnych w G , $V = \{S = X_1, X_2, \dots, X_N\}$, a p maksimum długości prawych stron produkcji. Weźmy słowo $\omega \in L(G)$ o długości $m > pN$. Zatem jego wyprowadzenie musi zawierać co najmniej $N + 1$ kroków

$$\begin{aligned} S = X_1 &\Rightarrow \alpha_1 X_{i_2} \eta_1 \Rightarrow \alpha_1 \alpha_2 X_{i_3} \eta_2 \eta_1 \Rightarrow \dots \Rightarrow \alpha_1 \dots \alpha_N X_{i_{N+1}} \eta_N \dots \eta_1 \\ &\Rightarrow \alpha_1 \dots \alpha_N \gamma \eta_N \dots \eta_1 = \omega. \end{aligned}$$

Stąd wniosek, że jeden z symboli X_i występuje w takim wyprowadzeniu co najmniej 2 razy

$$S \Rightarrow^* \alpha X_i \eta \Rightarrow^* \alpha \beta X_i \delta \eta \Rightarrow^* \alpha \beta \gamma \delta \eta,$$

a więc

$$X_i \Rightarrow^* \beta X_i \delta \quad \text{oraz} \quad X_i \Rightarrow^* \gamma.$$

Mamy zatem także

$$X_i \Rightarrow^* \beta^k \gamma \delta^k \quad \text{czyli} \quad \omega_k = \alpha \beta^k \gamma \delta^k \eta \in L(G) \quad \forall k \geq 0.$$

Lemat o pompowaniu dla jęz. liniowych

Niech L będzie nieskończonym językiem liniowym. Istnieje wówczas stała $m > 0$ (zależna jedynie od L) taka, że dowolne słowo $\omega \in L$ o długości co najmniej m można rozłożyć na części $\omega = \alpha \beta \gamma \delta \eta$ w taki sposób, że

- (i) $|\alpha \beta \delta \eta| \leq m$,
- (ii) $|\beta \delta| \geq 1$,

a przy tym wszystkie bez wyjątku słowa $\omega_k = \alpha \beta^k \gamma \delta^k \eta$ dla $k = 0, 1, 2, \dots$ są także elementami L .

Przykład. $L = \{\omega \in \{a, b\}^* : |\omega|_a = |\omega|_b\}$ nie jest liniowy. Załóżmy nie wprost, że istnieje generująca go liniowa gramatyka. Zgodnie z Lematem dla L istnieje stała m wyznaczająca długość słów ω , dla których zawsze możliwe jest pompowanie $\omega_k \in L$.

Weźmy $\omega = a^m b^{2m} a^m$. Zgodnie z lematem $\omega = \alpha\beta\gamma\delta\eta$, gdzie $|\alpha\beta| + |\delta\eta| \leq m$, a więc muszą być zbudowane wyłącznie z liter a , natomiast (ii) zapewnia, że pompowane części β i δ zawierają co najmniej jedną literę a . Część centralna słowa zbudowana z liter b ukryta jest w podciągu γ i nie podlega pompowaniu. Tak więc liczba liter a w słowach ω_k wzrasta, podczas gdy liczba liter b pozostaje niezmienną. Tym samym pompowanie tworzy słowa spoza L . Otrzymaliśmy sprzeczność z tezą lematu, a zatem założenie o liniowości języka L musiało być fałszywe.

Lemat o pompowaniu dla języków \mathcal{L}_2

Niech L będzie nieskończonym językiem bezkontekstowym. Istnieje wówczas stała $m > 0$ (zależna jedynie od L) taka, że dowolne słowo $\omega \in L$ o długości co najmniej m można rozłożyć na części $\omega = \alpha\beta\gamma\delta\eta$ w taki sposób, że

- (i) $|\beta\gamma\delta| \leq m$,
- (ii) $|\beta\delta| \geq 1$,

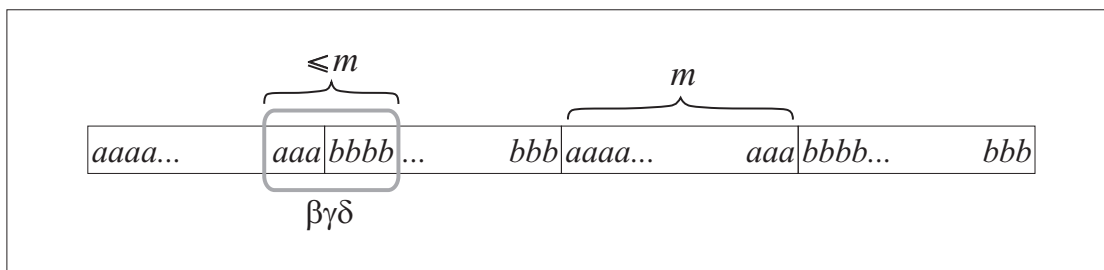
a przy tym wszystkie bez wyjątku słowa $\omega_k = \alpha\beta^k\gamma\delta^k\eta$ dla $k = 0, 1, 2, \dots$ są także elementami L .

Przykład. $L = \{a^n b^n c^n : n \geq 0\}$. Załóżmy, że jest to język bezkontekstowy. Istnieje zatem stała m gwarantująca, że słowo $\omega = a^m b^m c^m$ dzięki swojej długości może być podzielone na segmenty $\omega = \alpha\beta\gamma\delta\eta$ w sposób zgodny z (i) oraz (ii) tak, że zastosowana do niego operacja pompowania $\omega_k = \alpha\beta^k\gamma\delta^k\eta$ generuje inne słowa z L . Jednak warunek (ii) $|\beta\gamma\delta| \leq m$ dla naszego słowa oznacza, że segment ten może mieć tylko jedną z postaci

$$a^r, \quad a^s b^t, \quad b^r, \quad b^s c^t \quad \text{lub ostatecznie} \quad c^r,$$

gdzie $r \leq m$ lub $s + t \leq m$. Wówczas jednak pompowanie powieli tylko jedną lub dwie spośród liter a, b, c , produkując słowa spoza języka L . Stąd wniosek, że język ten nie może być bezkontekstowy.

Przykład. $L = \{\mu\mu : \mu \in \{a, b\}^*\}$. Niech $\omega = a^m b^m a^m b^m$.



Segment $\beta\gamma\delta$ ze względu na ograniczoną długość nie większą niż m może mieścić się w obrębie co najwyżej dwóch sąsiednich bloków liter a i b . W każdym z możliwych położen pompowanie niszczy symetrię słów języka L . Zatem nie może to być język bezkontekstowy.

Przykład. $L = \{a^{n^2} : n \geq 0\}$. Argument jest dokładnie *taki sam* jak w przypadku języków regularnych. Niech $L \in \mathcal{L}_2$ i $\omega = a^{m^2}$. Można więc je podzielić na 5 części

$$\omega = \alpha\beta\gamma\delta\eta = a^p a^q a^r a^s a^t, \quad p + q + r + s + t = m^2,$$

gdzie (i) $q + r + s \leq m$, (ii) $q + s \geq 1$. Na podstawie tezy lematu wszystkie słowa

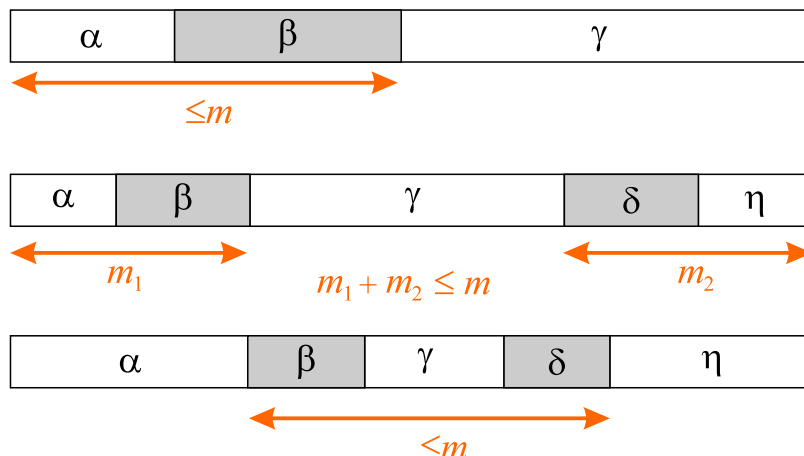
$$\omega_k = a^{p+(q+s)k+r+t}, \quad k = 0, 1, 2, \dots$$

są także elementami L . Jednak w L nie może być ani jednego słowa o pośredniej długości $m^2 < l < (m+1)^2$. Ale słowo $\omega_2 = a^{p+2q+r+2s+t} = a^{m^2+q+s}$ ma właśnie taką długość, bo $0 < q + s \leq m$, a więc $q + s < 2m + 1$ skąd:

$$m^2 < m^2 + q + s < m^2 + 2m + 1 = (m+1)^2.$$

Widać zatem, że założenie o bezkontekstowości L musi być fałszywe.

Lematy o pompowaniu – porównanie

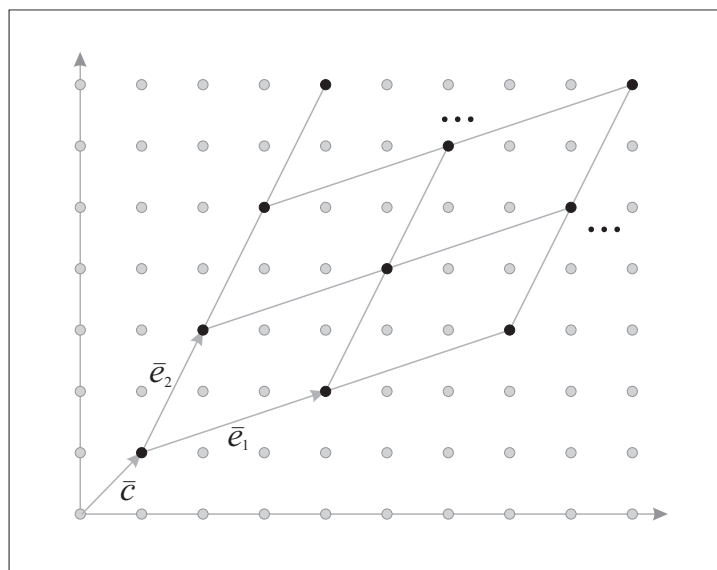


Twierdzenie Parikha

Podzbiór $K \subseteq \mathbb{N}^n$ nazywamy liniowym, jeśli istnieją wektory \bar{c} i $\bar{e}_1, \dots, \bar{e}_m \in \mathbb{N}^n$, $m \leq n$, dla których

$$K = \{ \bar{c} + l_1 \bar{e}_1 + l_2 \bar{e}_2 + \dots + l_m \bar{e}_m : l_1, l_2, \dots, l_m \in \mathbb{N} \}.$$

Zbiór K nazywamy semiliniowym, jeśli jest on sumą mnogościową skończonej liczby zbiorów liniowych.



Zbiór liniowy można uważać za wielowymiarowe uogólnienie ciągu arytmetycznego, $a_n = a_0 + nr$, gdzie \bar{c} pełni rolę wyrazu początkowego a_0 , a wektory \bar{e}_i są wielowymiarowymi przyrostami odpowiadającymi r .

Odwzorowaniem Parikha dla alfabetu $A = \{a_1, a_2, \dots, a_n\}$ nazywamy funkcję $\Psi : A^* \rightarrow \mathbb{N}^n$ określoną następująco:

$$\Psi(\omega) = (|\omega|_{a_1}, |\omega|_{a_2}, \dots, |\omega|_{a_n}).$$

Przykład: jeśli $A = \{a, b, c\}$, wówczas $\Psi(aacabb) = (3, 2, 1)$ itp.

Twierdzenie: *Jeśli język L jest bezkontekstowy, wówczas jego obrazem w odwzorowaniu Parikha $\Psi(L)$ jest zbiór semiliniowy.*

Zauważmy, że bardzo różne języki mogą mieć identyczną reprezentację w odwzorowaniu Parikha jako podzbiory \mathbb{N}^n . Na przykład $L = \{a^n b^n : n \geq 0\}$ i język poprawnie zagnieżdżonych nawiasów przekształcane są przez Ψ w zbiór $\{(n, n) : n \in \mathbb{N}\}$. Ma zatem sens następujące określenie: mówimy, że języki L i L' są literowo (permutacyjnie) równoważne jeśli $\Psi(L) = \Psi(L')$.

Lemat: *Każdy zbiór semiliniowy w \mathbb{N}^n jest obrazem pewnego języka regularnego nad n -literowym alfabetem.*

Uzasadnienie dla $n = 2$. Weźmy K rozpięte przez $\bar{c} = (c_1, c_2)$ oraz $\bar{e}_1 = (p, q)$ i $\bar{e}_2 = (r, s)$. Tworzymy na tej podstawie wyrażenie regularne

$$\rho = a^{c_1} b^{c_2} (a^p b^q)^* (a^r b^s)^*.$$

Łatwo sprawdzić, że $\Psi(L(\rho)) = K$.

Jeśli K jest semiliniowy, $K = K_1 \cup \dots \cup K_m$, utwórzmy podobnie wyrażenia regularne ρ_i dla K_i a następnie połączmy $\rho = \rho_1 + \dots + \rho_m$. $L(\rho)$ jest poszukiwanym językiem regularnym.

Wniosek: *Każdy język bezkontekstowy jest permutacyjnie równoważny z pewnym językiem regularnym.*

Stąd także: *Każdy język bezkontekstowy nad jednoliterowym alfabetem jest regularny.*

Uwaga: To że obraz Parikha jest zbiorem semiliniowym nie oznacza, że język jest bezkontekstowy, przykład: $L = \{a^n b^n c^n : n \geq 0\}$.

Wyniki dotyczące domkniętości i rozstrzygalności dla języków \mathcal{L}_2

Klasa \mathcal{L}_2 jest zamknięta ze względu na sumę mnogościową, konkatenację i $*$ Kleenego, nie jest natomiast zamknięta ze względu na przekrój i dopełnienie języków (a więc też nie ze względu na różnicę i różnicę symetryczną).

Istnieją przykłady języków bezkontekstowych, których przekroje i dopełnienia nadal są bezkontekstowe. Np. bezkontekstowe są języki

$$\{a^n b^n : n \geq 0\}^c, \quad \{\omega \overleftarrow{\omega} : \omega \in \{a, b\}^*\}^c$$

dla przekroju np.

$$\{a^n b^m : n \geq m\} \cap \{a^n b^m : n \leq m\} = \{a^n b^n : n \geq 0\}$$

Z kolei

$$\{a^n b^n c^m : n, m \geq 0\} \cap \{a^m b^n c^n : n, m \geq 0\} = \{a^n b^n c^n : n \geq 0\}$$

to przykład, gdy przekrój wyprowadza poza klasę \mathcal{L}_2 . Można też pokazać, że

$$L = \{\omega\omega : \omega \in \{a, b\}^*\}^c \quad \begin{array}{l} S \rightarrow A|B|AB|BA \\ A \rightarrow a|aAa|aAb|bAa|bAb \\ B \rightarrow b|aBa|aBb|bBa|bBb \end{array}$$

jest bezkontekstowy (gramatyka po prawej stronie). Wiemy jednak już, że L^c nie jest klasy \mathcal{L}_2 .

Oto jeden pozytywny rezultat:

Jeśli $L \in \mathcal{L}_2$ oraz $R \in \mathcal{L}_3$, wówczas $L \cap R$ jest bezkontekstowy.

Przykłady.

Pokażemy, że $L = \{a^n b^n : n \geq 0, n \neq 100\}$ jest bezkontekstowy. Weźmy jednoelementowy, a zatem regularny język $P = \{a^{100} b^{100}\}$. Dopełnienie $R = P^c$ też jest regularne. Ponieważ $L = \{a^n b^n : n \geq 0\} \cap R$, wnioskujemy, że jest on bezkontekstowy.

Pokażemy, że język $L = \{\omega : |\omega|_a = |\omega|_b = |\omega|_c\}$ nie jest bezkontekstowy. Wynika to natychmiast z obserwacji, że

$$L \cap L(a^* b^* c^*) = \{a^n b^n c^n : n \geq 0\}.$$

Zagadnienia rozstrzygalne w \mathcal{L}_2

1. Mając daną gramatykę G i słowo ω , rozstrzygnąć czy $\omega \in L(G)$.
(Postać normalna Chomsky'ego i algorytm CYK $O(n^3)$.)
2. Mając daną gramatykę G , rozstrzygnąć czy $L(G)$ jest pusty, skończony bądź nieskończony.
3. Mając dane 2 gramatyki G i H typu $LL(k)$, rozstrzygnąć czy $L(G) = L(H)$.

Zagadnienia nierozstrzygalne w \mathcal{L}_2

1. Rozstrzygnąć czy $L(G) = L(H)$.
2. Rozstrzygnąć czy $L(G) \subseteq L(H)$.
3. Zbadać czy $L(G) \cap L(H) = \emptyset$.
4. Zbadać czy $L(G) = A^*$.
5. Rozstrzygnąć czy $L(G)$ jest językiem regularnym (to jest, czy dla danej gramatyki bezkontekstowej istnieje równoważna gramatyka regularna).