

ROZDZIAŁ I

Podstawowe pojęcia i ogólne twierdzenia lingwistyki matematycznej

1. Wyrażenia, języki, gramatyki

Dowolny skończony niepusty zbiór V będziemy nazywać:

- 1) alfabetem, jeżeli będzie składać się on z liter,
 - 2) słownikiem, gdy jego elementami będą słowa,
- przyjmując jednak zarazem nazwę „alfabet” za uniwersalną.

Elementy V oznaczać będziemy małymi literami alfabetu łacińskiego (a, b, c, \dots) i nazywać je będziemy symbolami alfabetu (czy słownika) V .

Wyrażeniem nad V nazywać będziemy dowolny skończony ciąg symboli z V . Zbiór wszystkich wyrażeń nad V oznaczać będziemy przez V^* . Elementy z V^* oznaczać będziemy przez P, Q, R, \dots , przy czym dodatkowo obieramy symbol λ (czyt.: lambda) na oznaczenie pustego ciągu symboli z V . Zbiór wszystkich niepustych ciągów symboli nad V oznaczać będziemy przez V^+ (tj. $V^+ = V^* \setminus \{\lambda\}$).

Elementy V^* nazywamy:

- 1) słowami, gdy tworzone są z liter alfabetu V ,
- 2) daniami, gdy tworzone są ze słów słownika V .

Przykład 1.1.

Gdy $V = \{a, b\}$, to $V^* = \{\lambda, a, b, aa, ab, ba, bb, aaa, \dots\}$,
zaś $V^+ = \{a, b, aa, ab, ba, bb, aaa, \dots\}$. \square

Konkatenacją wyrażeń P i Q (gdzie $P = a_1 \dots a_n$, a $Q = b_1 \dots b_m$) nazywamy ich sklejenie, co oznaczamy i rozumiemy w następujący sposób:

$$P \cdot Q = PQ = a_1 \dots a_n b_1 \dots b_m.$$

λ jest oczywiście elementem neutralnym konkatenacji: $\lambda \cdot P = P \cdot \lambda = P$ (stąd $\lambda \lambda = \lambda$).

Bezpośrednio z definicji konkatenacji otrzymujemy, że jest ona łączna:

$P(QR) = (PQ)R$. Zwykle nie jest ona jednak przemienna (tj. zwykle $PQ \neq QP$).

Przykład 1.2.

Gdy $P = \text{kar}$, a $Q = \text{nawał}$, to $PQ = \text{karnawał}$, a $QP = \text{nawałkar}$. \square

Wprowadźmy kolejne definicje:

- 1) Długością wyrażenia $P \in V^*$ nazywamy liczbę $|P|$ określoną indukcyjnie:

$$|P| = \begin{cases} 0, & \text{gdy } P = \lambda \\ n, & \text{gdy } P = a_1 \dots a_n. \end{cases}$$

2) Przez P^n oznaczmy n -tą krotność wyrażenia P :

$$P^n = PP\dots P, \text{ gdzie } P \text{ występuje } n \text{ razy, przy czym } P^0 = \lambda.$$

3) Wyrażenie P nazywamy odcinkiem wyrażenia Q (co zapisujemy $P \subseteq Q$) wtedy, gdy istnieją pewne Q_1, Q_2 należące do V^* takie, że $Q = Q_1PQ_2$.

Zauważmy, że może być $Q_1 = \lambda$ lub $Q_2 = \lambda$.

4) Gdy w poprzedniej definicji $Q_1 = \lambda$, to mamy $Q = PQ_2$, a P nazywamy odcinkiem początkowym Q i fakt ten oznaczamy: $P \subseteq_e Q$.

5) Gdy z kolei w owej definicji $Q_2 = \lambda$, to mamy $Q = Q_1P$, a P nazywamy odcinkiem końcowym Q i fakt ten oznaczamy: $P \subseteq_f Q$.

6) Dla dowolnego słowa $P \in V^*$, przez P^{-1} określamy odbicie (zwierciadlane) wyrażenia P określone w następujący sposób: jeśli $P = a_1\dots a_n$, to $P^{-1} = a_n\dots a_1$.

Zauważmy, że zachodzi:

a) $(P^n)^{-1} = (P^{-1})^n$,

b) $(P^{-1})^{-1} = P$,

c) $(PQ)^{-1} = Q^{-1}P^{-1}$.

Zadanie 1.1: Wykaż zachodzenie powyższych zależności. \square

Przykład 1.2 - c.d.

Gdy $P=ma$, a $Q=ta$, to $PQ=mata$, $QP=tama$, a $PP=P^2=mama$ i $QQ=Q^2=tata$.

Nadto $P^{-1}=am$, a $(PQ)^{-1}=atam=Q^{-1}P^{-1}$. \square

Założmy, że zbiór symboli alfabetu V jest uporządkowany (tj. zadana jest kolejność jego elementów). Fakt, że litera $x \in V$ jest wcześniejszą od litery $y \in V$, zapisywać będziemy $x \prec y$. Tak rozumiana relacja \prec określa w zbiorze V porządek liniowy. Porządek ów możemy rozszerzyć na zbiór wszystkich słów w następujący sposób:

1) słowo puste poprzedza wszystkie pozostałe słowa,

2) litery alfabetu uporządkowane są według relacji \prec porządku liter,

3) wszystkie słowa zaczynające się od pierwszej litery alfabetu są wcześniejsze od słów zaczynających się od drugiej litery, itd. Jeżeli dwa słowa zaczynają się od tej samej litery alfabetu, to wcześniejsze jest to, którego druga litera jest wcześniejszą. Jeśli i drugie litery są takie same, to należy rozpatrywać trzecie litery, itd.... Jeżeli słowo P stanowi początek słowa Q , to słowo P jest wcześniejsze od słowa Q (chyba, że są one identyczne).

Porządek ten jest zgodny z porządkiem słów w słowniku. Dlatego też nosi on nazwę *porządku leksykograficznego*. Jest on bardzo wygodny dla skończonego zbioru słów (np. dla języka naturalnego i dlatego wykorzystuje się go m.in. w słownikach i encyklopediach). Jednak dla nieskończonego zbioru słów nad jakimś alfabetem jest on (jak to pokazuje poniższy przykład) dość dziwny i niepraktyczny.

Przykład 1.1 - c.d.

Dla $V = \{a, b\} - V^*$ z wypisanymi elementami w kolejności leksykograficznej wyglądałby następująco:

$$V^* = \{\lambda, a, aa, aaa, aaaa, \dots, aaab, aaaba, \dots, aab, \dots, ab, aba, abaa, \dots, b, ba, \dots\}. \quad \square$$

Widzimy, że jest on co pewien czas „poprzetykany” wielokropkiem zastępującym nieskończoną liczbę słów, a zatem jest on niewłaściwy dla maszyny (tj. komputera), która nie byłaby w stanie wygenerować w tym porządku niektórych słów, nawet dysponując nieskończonym czasem (w powyższym przykładzie, chociażby dowolnego słowa zawierającego chociażby jedno b, gdyż po słowie zawierającym n a, generowałyby słowo zawierające n+1 a, następnie n+2 a, ..., itd.). Dlatego też słowa nad alfabetem porządkuje się zazwyczaj w inny sposób, który nazywamy *porządkiem zgodnym z długością*, zgodnie z którym to słowa najpierw porządkuje się według długości, a potem dopiero słowa o równej długości porządkuje się leksykograficznie. Tak też zostały uporządkowane słowa w V^* w przykładzie 1.1 na początku tego rozdziału.

Zadanie 1.2: Pokaż, że:

- a) dla zwykłego porządku leksykograficznego,
 - b) dla porządku leksykograficznego zgodnego z długością,
- jeżeli x i y są takimi słowami, że $x \prec y$, a z jest dowolnym słowem, to $xz \prec yz$ oraz $zx \prec zy$. \square

Językiem (segmentowym, lub równoważnie ciągowym) L nad alfabetem V nazywamy dowolny zbiór słów utworzonych nad alfabetem V .

Język L jest więc pewnym podzbiorem zbioru V^* . Zauważmy ponadto, że tak sformułowana definicja nie precyzuje, które słowo jest, a które nie jest elementem języka. Podaje ona jedynie, że słowa w danym języku mogą być tworzone tylko nad danym alfabetem (tak więc np. słowo „vox” nie może być słowem języka polskiego, dopóki v i x nie staną się literami jego alfabetu). Ponieważ dla dowolnego V , zbiór V^* jest zawsze zbiorem nieskończonym, zatem z definicji języka wynika, że mogą być one tak zbiorami skończonymi, jak i nieskończonymi. Dany język możemy

dokładnie określić jedynie podając, z jakich słów się składa (przez ich wypisanie lub podanie ich kształtu).

Przykład 1.1 - c.d.

Nad alfabetem $V = \{a, b\}$ można określić m.in. następujące języki:

- 1) $L_1 = \emptyset$ (język pusty, tj. język w ogóle nie posiadający słów),
- 2) $L_2 = \{\lambda\}$ (język, którego jedynym elementem jest słowo puste),
- 3) $L_3 = \{\lambda, a, b\}$ (język ze słowami o długości co najwyżej 1),
- 4) $L_4 = \{a^i b^i : i = 0, 1, 2, \dots\} = \{\lambda, ab, aabb, aaabbb, \dots\}$,
- 5) $L_5 = \{PP^{-1} : P \in V^*\}$,
- 6) $L_6 = \{a^{n^2} : n = 0, 1, 2, \dots\}$,
- 7) Niech $N_a(P)$ - oznacza liczbę występowania symbolu a w P ,
a $N_b(P)$ - liczbę występowania symbolu b w P .

Wówczas $L_7 = \{P : P \in \{a, b\}^+ \wedge N_a(P) = N_b(P)\}$ (np. $ababba \in L_7$). \square

Zadanie 1.3: Opisz języki $L_5 - L_7$. \square

Dla celów następczej definicji przyjmijmy oznaczenie, że gdy $A, B \subset V^*$, to

$$AB = \{P_1 P_2 : P_1 \in A \wedge P_2 \in B\}.$$

Gramatyką (generatywną) nazywamy uporządkowaną czwórkę $G = \langle V_N, V_T, S, F \rangle$, w której:

- 1) V_N i V_T są niepustymi, skończonymi i rozłącznymi alfabetami,
- 2) $S \in V_N$,
- 3) $F \subseteq (V_N \cup V_T)^* V_N (V_N \cup V_T)^* \times (V_N \cup V_T)^*$.

V_N - to alfabet nieterminalny (pomocniczy), złożony ze zmiennych syntaktycznych (lub równoważnie: metajęzykowych albo pomocniczych).

V_T - to alfabet symboli terminalnych (przedmiotowych), albo krócej - alfabet terminalny (końcowy).

S - to symbol początkowy.

F - to zbiór par (P, Q) , takich że w P występuje zawsze przynajmniej jeden symbol nieterminalny, a Q jest dowolnym wyrażeniem o elementach należących do $V_N \cup V_T$.

Elementy F nazywamy regułami produkcji (lub przepisywania).

Zamiast pisać (P, Q) , piszemy zwykle $P \rightarrow Q$.

P nazywamy tu poprzednikiem, a Q następnikiem produkcji.

Zasada stosowania reguł produkcji jest następująca: jeśli mamy wyrażenie, którego częścią jest wyrażenie P , a w gramatyce, którą się posługujemy, jest reguła produkcji $P \rightarrow Q$, wówczas możemy wyrażenie to zastąpić nowym wyrażeniem, w którym na miejscu P jest Q . Zobaczmy to na poniższym przykładzie:

Przykład 1.3.

Niech $G = \langle V_N, V_T, S, F \rangle$, gdzie

$V_N = \{ \langle \text{zdanie} \rangle, \langle \text{podmiot} \rangle, \langle \text{orzeczenie} \rangle \}$ (jest to zbiór symboli nieterminalnych, czyli pomocniczych, zwanych tu zmiennymi metajęzykowymi),

$V_T = \{ a, \acute{a}, b, \dots, z, \acute{z}, \acute{z} \}$ (alfabet języka polskiego),

$S = \langle \text{zdanie} \rangle$ (symbol początkowy, określający końcowy wynik działania gramatyki),

$F = \{ \langle \text{zdanie} \rangle \rightarrow \langle \text{podmiot} \rangle \langle \text{orzeczenie} \rangle, \langle \text{podmiot} \rangle \rightarrow \text{Jaś}, \langle \text{orzeczenie} \rangle \rightarrow \acute{s}pi \}$ (zbiór reguł produkcji).

W gramatyce tej wyrażenie „Jaś śpi” otrzymujemy następująco:

$\langle \text{zdanie} \rangle \rightarrow \langle \text{podmiot} \rangle \langle \text{orzeczenie} \rangle \rightarrow \text{Jaś} \langle \text{orzeczenie} \rangle \rightarrow \text{Jaś} \acute{s}pi \quad \square$

Mówimy, że gramatyka G przekształca bezpośrednio wyrażenie $P \in (V_N \cup V_T)^*$ (co czytamy: P utworzone nad sumą alfabetów terminalnego i nieterminalnego) w wyrażenie $Q \in (V_N \cup V_T)^*$, co oznaczamy: $P \xrightarrow{G} Q$, witw, gdy stosując jeden raz jedną z reguł produkcji gramatyki G , z wyrażenia P otrzymujemy wyrażenie Q .

Uogólnijmy pojęcie przekształcania bezpośredniego.

Mówimy, że gramatyka G przekształca słowo $P \in (V_N \cup V_T)^*$ w słowo $Q \in (V_N \cup V_T)^*$ (co oznaczamy: $P \xrightarrow{*G} Q$) witw, gdy istnieje ciąg wyrażen

$P = P_0, P_1, \dots, P_n = Q$, taki że $P = P_0 \xrightarrow{G} P_1 \xrightarrow{G} P_2 \xrightarrow{G} \dots \xrightarrow{G} P_n = Q$.

Taki ciąg P_0, P_1, \dots, P_n nazywamy wyprowadzeniem (derywacją) wyrażenia Q z P .

Przykład 1.4.

Niech $V = \{ a, b, c, +, \cdot, (,) \}$. Skonstruujemy gramatykę złożoną ze wszystkich poprawnie zbudowanych wyrażen algebraicznych z elementów zbioru V :

$V_N = \{ S \},$

$V_T = V,$

$F = \{ S \rightarrow S + S, S \rightarrow S \cdot S, S \rightarrow (S + S), S \rightarrow (S \cdot S), S \rightarrow a, S \rightarrow b, S \rightarrow c \}.$

W gramatyce tej wyrażenie $((a \cdot b) + c)$ otrzymujemy następująco:

$$S \xrightarrow{G} (S+S) \xrightarrow{G} ((S \cdot S)+S) \xrightarrow{G} ((a \cdot S)+S) \xrightarrow{G} ((a \cdot b)+S) \xrightarrow{G} ((a \cdot b)+c).$$

Tak więc gramatyka ta przekształca bezpośrednio np. wyrażenie $((a \cdot S)+S)$ w wyrażenie $((a \cdot b)+S)$. Podobnie możemy powiedzieć, że przekształca ona wyrażenie $(S+S)$ w wyrażenie $((a \cdot b)+S)$. \square

Zauważmy, że:

- 1) zmienne alfabetu nieterminalnego są rzeczywiście jedynie zmiennymi pomocniczymi. W powyższym przykładzie służą one ponadto do wyodrębnienia w wyrażeniu grup (wraz z regułami produkcji są więc tu odpowiedzialne za jego syntaktykę). Podobnie ma się sprawa w przypadku przykładu 1.3;
- 2) gdy wiadomo jest, z jakiej gramatyki korzystamy, to przy kolejnych krokach derywacji nie musimy pisać pod strzałką „G”, specyfikującego tę gramatykę.

Zadanie 1.4: W gramatyce z przykładu 1.4 wyprowadź wyrażenie

$$(((a+b) \cdot c + (b+a) \cdot a) \cdot b). \quad \square$$

Zadanie 1.5: Udowodnij następujące własności wynikania bezpośredniego:

- a) $v \xrightarrow{*} v$ dla dowolnego $v \in V^*$,
- b) jeżeli $v \xrightarrow{*} u$ i $u \xrightarrow{*} w$, to $v \xrightarrow{*} w$ dla dowolnych $u, v, w \in V^*$,
- c) jeżeli $v \xrightarrow{*} u$, to $vy \xrightarrow{*} uy$ i $yv \xrightarrow{*} yu$ dla dowolnych $u, v, y \in V^*$. \square

Dla P i Q takich, że $P \xrightarrow{*} Q$, do Q z P często można dojść na wiele sposobów,

i to w dodatku mogą się różnić one swą długością. W takiej sytuacji za długość derywacji będziemy przyjmować długość najkrótszej derywacji prowadzącej z P do Q .

Językiem $L(G)$ generowanym przez gramatykę generatywną G nazywamy zbiór

$$L(G) = \{P \in V_T^* : S \xrightarrow{*}_G P\}.$$

Jest to więc zbiór tych słów utworzonych nad alfabetem terminalnym gramatyki G , które są w niej wyprowadzalne z jej symbolu początkowego S .

Tak określony język generowany przez gramatykę jest oczywiście językiem (tj. spełnia definicję języka). Składa się on bowiem ze słów utworzonych nad pewnym alfabetem (tu: V_T); jest on podzbiorem (wyznaczonym przez wyprowadzalność z S) zbioru wszystkich słów z V_T^* .

Mówimy, że dwie gramatyki G_1 i G_2 są (słabo) równoważne w tw, gdy języki przez nie generowane są identyczne (tj. $L(G_1) = L(G_2)$).

Z kolei mówimy, że dwie gramatyki są mocno równoważne, gdy są one słabo równoważne i identyczne są odpowiednie drzewa derywacji w tych gramatykach (z którymi szczegółowiej zaznajomimy się w dalszej części kursu).

Poniżej podano dwa przykłady równości języków (generowanego przez gramatykę i zdefiniowanego przez podanie konstytuujących go słów).

Przykład 1.5.

Niech $G = \langle V_N, V_T, S, F \rangle$, $V_N = \{S, A, B\}$, $V_T = \{a, b\}$,

$F = \{S \rightarrow aB, B \rightarrow aBB, B \rightarrow b, B \rightarrow bS, S \rightarrow bA, A \rightarrow bAA, A \rightarrow a, A \rightarrow aS\}$.

Wówczas $L(G) = L_7$ (gdzie L_7 jest językiem zdefiniowanym w przykładzie 1.1 na str. 9). \square

Przykład 1.6.

Niech $G = \langle \{S, X, Y\}, \{a, b, c\}, S, F \rangle$, gdzie

$F = \{S \rightarrow abc, S \rightarrow aXbc, Xb \rightarrow bX, Xc \rightarrow Ybcc, bY \rightarrow Yb, aY \rightarrow aaX, aY \rightarrow aa\}$.

Wówczas $L(G) = \{a^n b^n c^n : n \geq 1\}$. \square

Zauważmy na koniec, że wyprowadzanie wyrażenia z S może się zatrzymać wyłącznie w dwóch przypadkach:

1) gdy uzyskamy wyrażenie składające się wyłącznie z elementów z V_T

$(S \rightarrow \dots \rightarrow P \in V_T^*)$,

2) gdy dostaniemy wyrażenie $P \in (V_N \cup V_T)^*$ takie, że nie ma reguł produkcji, które to byśmy mogli do niego zastosować.

(1) jest oczywiście szczególnym przypadkiem 2)).

Zadanie 1.6: Niech $G = \langle V_N, V_T, S, F \rangle$ i $G' = \langle V_N, V_T, S, F' \rangle$, gdzie $F' \subset F$.

Pokaż, że $L(G') \subseteq L(G)$. \square

Zadanie 1.7: Udowodnij, że jeśli $x \xrightarrow{*} y$ i $x = x_1 x_2 \dots x_k$, to istnieją słowa

y_1, y_2, \dots, y_k , takie że $x_1 \xrightarrow{*} y_1, \dots, x_k \xrightarrow{*} y_k$, a $y = y_1 y_2 \dots y_k$. \square

Zadanie 1.8: Nad alfabetem języka polskiego zdefiniuj gramatykę generującą tylko i wyłącznie wszystkie słowa będące imionami studentek i studentów z twojej grupy. \square

Zadanie 1.9: Niech symbol „_” oznacza spację. Skonstruuj gramatykę generującą tylko i wyłącznie wszystkie słowa postaci „Mam_Q.”, gdzie zamiast Q stoi dowolny niepusty ciąg liter utworzony nad alfabetem polskim. \square

2. Hierarchia Chomsky’ego

Amerykański lingwista N. Chomsky sklasyfikował gramatyki ze względu na kształt dopuszczalnych w nich reguł produkcji. Zgodnie z ową klasyfikacją mówimy, że gramatyka generatywna $G = \langle V_N, V_T, S, F \rangle$ jest typu (lub równoważnie klasy):

- 1) 0 (lub gramatyką struktur frazowych), gdy wszystkie jej reguły produkcji mają postać $P \rightarrow Q$, gdzie $P \in (V_N \cup V_T)^* V_N (V_N \cup V_T)^*$, zaś $Q \in (V_N \cup V_T)^*$ (tj., w stosunku do definicji gramatyki, brak jest w niej jakichkolwiek ograniczeń na kształt reguł produkcji);
- 2) 1 (lub gramatyką kontekstową bądź czułą na kontekst, a po angielsku: context-sensitive grammar, lub krócej CS-grammar), gdy wszystkie jej reguły produkcji mają postać $Q_1 A Q_2 \rightarrow Q_1 P Q_2$, gdzie $Q_1, Q_2 \in (V_N \cup V_T)^*$, $A \in V_N$, zaś $P \in (V_N \cup V_T)^* \setminus \{\lambda\}$, bądź też $S \rightarrow \lambda$, ale wówczas S nie występuje po prawej stronie w żadnej z reguł produkcji;
- 3) 2 (lub bezkontekstową, a po angielsku: context-free grammar, lub krócej CF-grammar), gdy wszystkie jej reguły produkcji mają postać $A \rightarrow P$, gdzie $A \in V_N$, zaś $P \in (V_N \cup V_T)^*$ (tj. są postaci $\lambda A \lambda \rightarrow \lambda P \lambda$);
- 4) 3 (lub równoważnie gramatyką: regularną, prawoliniową, prawostronnie liniową, automatową), gdy wszystkie jej reguły produkcji mają postać $A \rightarrow P$ lub $A \rightarrow PB$, gdzie $A, B \in V_N$, zaś $P \in V_T^*$.

Jak się przekonamy, każda gramatyka klasy i ($i=0,1,2,3$) jest równocześnie gramatyką klasy j , dla $0 \leq j \leq i$. Mamy więc rzeczywiście do czynienia z hierarchią (wzajemnym zawieraniem się zbiorów).

Język L nazywamy językiem typu i , gdy jest on generowany przez pewną gramatykę typu i (tj. $L = L(G)$, gdzie G jest gramatyką typu i).

Klasę języków typu i oznaczamy \mathcal{L}_i .

Z powyższych definicji widzimy, że $\mathcal{L}_3 \subseteq \mathcal{L}_2 \subseteq \mathcal{L}_0$ oraz $\mathcal{L}_1 \subseteq \mathcal{L}_0$.

W dalszym ciągu pokażemy, że zachodzi następujące właściwe zawieranie się klas języków: $\mathcal{L}_3 \subset \mathcal{L}_2 \subset \mathcal{L}_1 \subseteq \mathcal{L}_0$ (tj., że ponadto $\mathcal{L}_2 \subseteq \mathcal{L}_1$ oraz że między \mathcal{L}_1 , \mathcal{L}_2 i \mathcal{L}_3 zachodzą inkluzje właściwe). W rzeczywistości zachodzi nawet: $\mathcal{L}_3 \subset \mathcal{L}_2 \subset \mathcal{L}_1 \subset \mathcal{L}_0$, tj. dodatkowo \mathcal{L}_1 właściwie zawiera się w \mathcal{L}_0 , ale w ramach naszego kursu faktu tego formalnie nie wykażemy.

3. Operacje na językach. Zamkniętość klas języków ze względu na operacje regularne

Ponieważ (zgodnie z definicją) języki są zbiorami, zatem możemy wykonywać na nich przynależne im operacje. Dla dowolnych języków L_1, L_2 (o słowach z V^*), mamy więc:

$L_1 \cup L_2 = \{P: P \in L_1 \vee P \in L_2\}$ (suma języków);

$L_1 \cap L_2 = \{P: P \in L_1 \wedge P \in L_2\}$ (iloczyn języków);

$L_1 \setminus L_2 = \{P: P \in L_1 \wedge P \notin L_2\}$ (różnica języków);

$\overline{L_1} = V^* \setminus L_1$ (dopełnienie języka);

Zadanie 1.10: Wykaż, że dla dowolnego języka L :

$$L \cup \emptyset = \emptyset \cup L = L, \text{ a } L \cup L = L. \quad \square$$

Zadanie 1.11: Niech $V = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$.

- Co to jest V^* ? Jakie znaczenia mają słowa z V^* ? Czy są to liczby?
- Określ tak język L nad alfabetem V , by był on złożony (tylko i wyłącznie) ze wszystkich liczb naturalnych. Zdefiniuj język \overline{L} .
- Jakie własności miałby ten język L , jeśli dodatkowo spełniałby warunek: $y \in L \Leftrightarrow y^{-1} \in L$. Co powiesz o analogicznym przypadku dla \overline{L} ? \square

$L_1 L_2 = \{P_1 P_2 : P_1 \in L_1 \wedge P_2 \in L_2\}$ (konkatenacja języków),

przy czym, ponieważ możemy tworzyć konkatenację (i to wielokrotną) tylko jednego języka (postaci: LL, LLL, \dots), więc definiujemy też:

$$L^i = \begin{cases} \{\lambda\}, & \text{gdy } i = 0 \\ L^k L, & \text{gdy } i = k + 1; \end{cases}$$

Zadanie 1.12: Wykaż, że $\emptyset L = L \emptyset = \emptyset$, a z kolei $\{\lambda\}L = L\{\lambda\} = L$

(WSKAZÓWKA: „ $\emptyset L$ ” oznacza język, w którym nic nie może stać na pierwszym miejscu). \square

Zadanie 1.13: Niech K , L i M będą trzema dowolnymi językami. Udowodnij, że $(LM)K = L(MK)$, tzn. że operacja konkatencji języka jest łączna. \square

$$L^* = \bigcup_{i \geq 0} L^i \text{ (domknięcie Kleene'go języka } L\text{);}$$

$$L^+ = \bigcup_{i > 0} L^i \text{ (niepełne domknięcie Kleene'go języka } L\text{);}$$

Zadanie 1.14: Udowodnij, że zachodzi:

$$\text{a) } L^* = L^+ \Leftrightarrow \lambda \in L,$$

$$\text{b) } L^+ = L^* \setminus \{\lambda\} \Leftrightarrow \lambda \notin L. \quad \square$$

Dla dowolnego języka L : $L^{-1} = \{P^{-1} : P \in L\}$ (odbicie lustrzane języka L).

Zadanie 1.15: Wykaż, że:

$$\text{a) } (L^{-1})^i = (L^i)^{-1},$$

$$\text{b) } (L^{-1})^{-1} = L. \quad \square$$

Operacjami regularnymi na językach nazywa się sumę języków, konkatencję języków i domknięcie Kleene'go języka.

Twierdzenie 1.1.

Każda z klas języków typu i (dla $i = 0, 1, 2, 3$) jest zamknięta ze względu na operacje regularne.

Ponieważ **dowód** tego twierdzenia jest dość pracochłonny (zajmuje kilka stron) - więc go pomijamy, tym bardziej, że twierdzenie to w pewnej modyfikacji i zawężeniu zostanie podane i udowodnione w trzecim rozdziale niniejszego podręcznika (tw. 3.3.). \square

4. I twierdzenie o postaci normalnej

Aby móc udowodnić tytułowe twierdzenie niniejszego paragrafu, musimy najpierw zdefiniować pojęcie homomorfizmu.

Niech mianowicie V_1 i V_2 będą dwoma alfabetami.

Odwzorowanie $h: V_1^* \rightarrow V_2^*$ nazywamy **h o m o m o r f i z m e m**, gdy:

$$\text{a) } \forall P \in V_1^* \exists! W \in V_2^* : W = h(P),$$

$$\text{b) } \forall P, Q \in V_1^* \quad h(PQ) = h(P)h(Q)$$

(symbol „ $\exists!$ a” oznacza: istnieje dokładnie jedno a).

Tak więc jest to odwzorowanie jednoznaczne, zawsze istnieje, a ponadto zachowuje ono konkatencję wyrażeń. W związku z tym, dowolny homomorfizm wystarczy określić jedynie dla alfabetu, bo gdy $P = a_1 \dots a_n$, to $h(P) = h(a_1) \dots h(a_n)$.

Oczywiście $h(\lambda)=\lambda$, bo: gdy $P \neq \lambda$, to $h(\lambda P) = h(\lambda)h(P) \rightarrow h(\lambda) = \lambda$.

$$\begin{array}{c} \parallel \\ h(P) \end{array} \quad \left. \begin{array}{c} \\ \end{array} \right\} \rightarrow h(\lambda) = \lambda$$

Homomorfizm nazywamy λ - wolnym, gdy $\forall P \neq \lambda \quad h(P) \neq \lambda$.

Gdy $L \subseteq V_1^*$, a $h: V_1^* \rightarrow V_2^*$, to homomorficzny obraz języka L definiujemy w następujący sposób: $h(L) = \{h(P): P \in L \subseteq V_1^*\}$.

Obecnie możemy przejść do podania drugiego już z naczelných twierdzeń lingwistyki matematycznej.

Twierdzenie 1.2 (I twierdzenie o postaci normalnej)

Dla każdej gramatyki $G_i = \langle V_N, V_T, S, F \rangle$ ($i=0, 1, 2, 3$), istnieje równoważna jej gramatyka $G_i' = \langle V_N', V_T, S, F' \rangle$ tego samego typu, taka że symbole terminalne nie występują po lewej stronie reguł produkcji w F' .

Dowód.

Dla gramatyk typu 2 i 3 nie ma czego dowodzić, bo wszystkie reguły produkcji mają w nich postać:

- 1) w gramatykach typu 2: $A \rightarrow P$, gdzie $A \in V_N$, a $P \in (V_N \cup V_T)^*$,
 - 2) w gramatykach typu 3: $A \rightarrow P$ lub $A \rightarrow PB$, gdzie $A, B \in V_N$, a $P \in V_T^*$,
- tj. po ich lewych stronach nie występują symbole terminalne.

Wystarczy więc udowodnić twierdzenie w przypadku, gdy $i=0$ i gdy $i=1$ (gdzie symbole terminalne po lewych stronach reguł produkcji mogą wystąpić).

W naszej gramatyce $G = \langle V_N, V_T, S, F \rangle$, niech $V_T = \{a_1, \dots, a_k\}$.

Dla każdego $a_i \in V_T$, obieramy nowy symbol nieterminalny $A_i \notin V_N$. Wówczas definiujemy $V_N' = V_N \cup \{A_1, \dots, A_k\}$, a następnie tworzymy F' w następujący sposób:

- 1) każdą regułę produkcji $\lceil P \rightarrow Q \rceil \in F$ zastępujemy nową regułą, w której wszystkie symbole terminalne występujące w starej regule zastąpione zostają odpowiadającymi im symbolami nieterminalnymi (tak więc np. regułę $\lceil Aa_1B \rightarrow a_1a_2a_3 \rceil \in F$, zastępujemy regułą $\lceil AA_1B \rightarrow A_1A_2A_3 \rceil \in F'$),
- 2) dla każdego $1 \leq i \leq k$, do F' włączamy reguły postaci $A_i \rightarrow a_i$.

W ten sposób utworzyliśmy gramatykę $G' = \langle V_N', V_T, S, F' \rangle$, w której rzeczywiście symbole terminalne nie występują po lewej stronie reguł produkcji.

Pozostaje nam jeszcze jedynie pokazać, że $L(G) = L(G')$.

- a) $L(G) \subseteq L(G')$, bo jeśli $a_{i_1} \dots a_{i_n} \in L(G)$, to $S \xrightarrow{*}_G A_{i_1} \dots A_{i_n} \xrightarrow{*}_G a_{i_1} \dots a_{i_n} \in L(G')$.

Najpierw stosowaliśmy tu reguły typu 1), odpowiednie do reguł gramatyki G , gwarantujące nam uzyskanie takich samych słów, jednak zbudowanych jedynie z symboli nieterminalnych, a następnie reguły typu 2), zamieniające poszczególne nieterminały na odpowiadające im terminały. W wyniku tych operacji, w gramatyce G' , uzyskiwaliśmy poszczególne słowa języka $L(G)$.

b) Pokażemy, że $L(G) \supseteq L(G')$.

W tym celu określimy homomorfizm $h: (V_{N'} \cup V_T)^* \rightarrow (V_N \cup V_T)^*$ w następujący sposób:

$$1) \forall 1 \leq i \leq k \quad h(A_i) = a_i,$$

$$2) h(x) = x \text{ dla pozostałych elementów z alfabetu } V_{N'} \cup V_T.$$

Tak określona funkcja h działa na symbolach alfabetu w odwrotną stronę niż w przypadku funkcji zmiany liter w regułach produkcji F na reguły produkcji F' .

Mamy wówczas, że gdy $P, Q \in (V_{N'} \cup V_T)^*$, a $P \xrightarrow[G']{*} Q$, to $h(P) \xrightarrow[G]{*} h(Q)$ (wniosek ten

otrzymujemy bezpośrednio z definicji homomorfizmu h , rozpatrując oddzielnie stosowanie obydwu typów reguł produkcji z F'). W szczególności, obierając S za P , otrzymujemy że:

gdy $S \xrightarrow[G']{*} Q$ (gdzie $Q \in V_T^*$), to $S = h(S) \xrightarrow[G]{*} h(Q) = Q$ (gdzie równości

uzyskaliśmy na mocy 2)). Zatem $L(G) \subseteq L(G')$, co już kończy dowód. \square