

Környezetfüggetlen nyelvek

Kiegészítő anyag az Algoritmuselmélet tárgyhoz VI.

(a Rónyai–Ivanyos–Szabó: Algoritmusok könyv mellé)

Friedl Katalin
BME SZIT
friedl@cs.bme.hu

2017. február 20.

A reguláris nyelveket véges automatákkal vagy reguláris kifejezésekkel adtuk meg. Van további lehetőség is, egy ilyen a *formális nyelvtan* vagy röviden *nyelvtan*.

A formális nyelvtanok nem egészen olyanok, mint például a magyar nyelvtan. Bár eredetileg a természetes nyelvek szabályainak leírására készültek, de könnyebben használhatóak mesterséges nyelvek, például programozási nyelvek pontos megadására, mint egy beszélt nyelv helyes mondatainak tökéletes leírására.

A formális nyelvtanok lényegében átírási szabályokat adnak meg, amelyekkel egy kezdő szimbólumból kiindulva szavakat tudunk előállítani.

Ezeknek itt csak egy speciális, talán a leggyakrabban használt fajtájával foglalkozunk.

1. Definíció. *Egy környezetfüggetlen nyelvtan (röviden CF nyelvtan) alatt egy olyan $G = (V, \Sigma, S, P)$ rendszert értünk, ahol*

- V egy véges, nem üres halmaz, a változók (vagy nemterminálisok) halmaza,
- Σ egy ábécé, amire $V \cap \Sigma = \emptyset$, a karakterek (vagy terminálisok) halmaza
- $S \in V$ a kezdő változó,
- P egy véges halmaz, az ún. levezetési (vagy produkciós, illetve átírási) szabályok halmaza.

P elemei $A \rightarrow \alpha$ alakúak, ahol $A \in V$ egy változó, $\alpha \in (V \cup \Sigma)^*$ egy változókból és a Σ elemeiből álló tetszőleges, véges hosszú sorozat.

1. Példa. Legyen $V = \{A\}$, $\Sigma = \{a, b\}$, a kezdő változó természetesen az A , a levezetési szabályok halmaza pedig

$$P = \{A \rightarrow aAb, A \rightarrow \varepsilon\}$$

A nyelvtanok megadásánál sokszor nem írjuk ki az összes paramétert, csak a levezetési szabályokat soroljuk fel. Ha mást nem mondunk, akkor a szabályokban szereplő kisbetűk a Σ elemei, a nagybetűk a változók, az első szabály bal oldala a kezdő változó. Továbbá azok a szabályok, melyeknek a bal oldalán ugyanaz áll, összevonhatóak, függőleges vonallal elválasztva a különböző jobb oldalakat, Ebben a formában az előző nyelvtan így néz ki:

$$A \rightarrow aAb \mid \varepsilon$$

2. Definíció. Egy $G = (V, \Sigma, S, P)$ nyelvtannál levezetés alatt egy olyan

$$\gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_n$$

véges hosszú sorozatot értünk ($n \geq 0$), melyben $\gamma_0 = S$, továbbá $\gamma_i \in (V \cup \Sigma)^*$, és mindegyik γ_{i+1} megkapható γ_i -ből egy levezetési szabály alkalmazásával. Ez azt jelenti, hogy minden $0 \leq i < n$ esetén γ_i felírható $\gamma_i = \delta_1 A \delta_2$ alakban, ahol $\delta_1, \delta_2 \in (V \cup \Sigma)^*$ és $A \in V$ úgy, hogy $\gamma_{i+1} = \delta_1 \alpha \delta_2$ ahol $A \rightarrow \alpha$ egy P -hez tartozó levezetési szabály.

2. Példa. Az előző nyelvtan esetén egy levezetés pl. az alábbi

$$A \Rightarrow aAb \Rightarrow aaAbb \Rightarrow aaaAbbb \Rightarrow aaa\varepsilonbbb$$

Az így levezetett szó az $aaabbb$. Könnyű látni, hogy ebből a nyelvtanból pontosan azok a szavak vezethetők le, amelyek $a^n b^n$ alakúak ($n \geq 0$).

A levezetések közül azok lesznek számunkra érdekesek, amelyekben a kezdő változóból indulva végül egy olyan sorozathoz jutunk, amiben már nincsenek változók.

3. Definíció. A $G = (V, \Sigma, S, P)$ nyelvtan által generált $L(G)$ nyelv azokból a $w \in \Sigma^*$ szavakból áll, melyekhez valamilyen $n \geq 0$ számra van olyan $S \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_n$ levezetés, amiben $\gamma_n = w$

1. Megjegyzés. Vegyük észre, hogy ha egy levezetés során eljutunk egy $w \in \Sigma^*$ szóhoz, akkor a levezetés tovább már biztos nem folytatható, mivel w nem tartalmaz változót, ilyenkor már egyetlen szabály sem alkalmazható.

3. Példa. Tekintsük az alábbi nyelvtant!

$$S \rightarrow aS \mid bS \mid a \mid b$$

Itt tehát az egyetlen változó az S , az ábécé az $\{a, b\}$. A nyelvtanból levezethető pl. az ab szó:

$$S \Rightarrow aS \Rightarrow ab$$

Az is látszik, hogy a nyelvtan által generált nyelv $\{a, b\}^*$ -ből az összes nem üres szót tartalmazza, azaz $L = \{a, b\}^* \setminus \{\varepsilon\}$. Ez azért igaz, mert egy tetszőleges, legalább 1 hosszú $w \in \{a, b\}^*$ szónak előlről kezdve egymás után tudjuk generálni a karaktereit: amíg nem az utolsó karakterről van szó, addig az vagy az első vagy a második szabállyal, az utolsó karakter pedig a 3. vagy a 4. szabállyal állítható elő.

1. Feladat. Mely szavakból áll az alábbi nyelvtan által generált nyelv?

$$\begin{aligned} S &\rightarrow aT a \mid bT b \mid a \mid b \\ T &\rightarrow aT \mid bT \mid \varepsilon \end{aligned}$$

Megoldás: Világos, hogy S -ből csak olyan szó vezethető le, ami nem üres, továbbá az első és az utolsó karaktere megegyezik. Megmutatjuk, hogy a generált nyelv az összes ilyen szóból áll.

Ehhez vegyük észre, hogy T -ből minden a és b betűkből álló sorozat előállítható az előző példához hasonlóan. Az S szabályai lehetővé teszik, hogy a szó első és utolsó karakterét generáljuk. Ha 1-nél hosszabb szót akarunk, akkor ezek közé a T -ből tetszőleges karaktersorozatot előállíthatunk. \square

2. Feladat. Legyen $\Sigma = \{a, b\}$ és L álljon azokból a szavakból, melyekben az a betűk száma megegyezik a b betűk számával. Adjunk olyan G nyelvtant, amire $L(G) = L$.

Megoldás: Egy lehetséges megoldás: $S \rightarrow aSbS \mid bSaS \mid \varepsilon$

Ahhoz, hogy ez valóban jó nyelvtan, először is vegyük észre, hogy minden esetben, amikor valamelyik szabályt alkalmazzuk, ugyanannyi a -t generálunk, mint ahány b -t ezért $L(G) \subseteq L$.

Azt kell még megmutatni, hogy minden $w \in L$ szó levezethető a nyelvtanból. Ezt a w hossza szerinti indukcióval látjuk be. Nyilván ez igaz a 0 hosszú $w = \varepsilon$ szóra. Egy hosszú szó nincs az L nyelvben. A kettő hosszú szavakra is könnyű látni, mert vagy az első vagy a második szabály egyszeri alkalmazása után a harmadik szabályt kétszer használva megkaphatjuk a w szót.

Tegyük fel, hogy L legfeljebb k hosszú szavairól már tudjuk, hogy levezethetők és legyen $|w| = k + 1$. Több eset lehetséges: amennyiben $w = aw'b$, akkor $w' \in L$ és ebben az esetben az $S \Rightarrow aSbS$ kezdés után S első előfordulásából, az indukciós feltevés szerint w' levezethető, miután a második S betűre az $S \rightarrow \varepsilon$ szabályt alkalmazva megkapjuk a w szót.

Hasonlóan járhatunk el, amennyiben $w = bw'a$.

Ha viszont w első és utolsó betűje megegyezik, és ez a betű mondjuk a , akkor vegyük w -nek egy legrövidebb kezdőszelétét, amiben ugyanannyi a van mint b , legyen ez x és $w = xy$. Ekkor egyrészt $x, y \in L$, másrészt x szükség szerűen b -re végződik, azaz $x = azb$ alakú, ahol $z \in L$. Egy ilyen w -re jó levezetést kapunk, ha az $S \Rightarrow aSbS$ lépés után az első S -ből az x -et, a másodikból az y -t vezetjük le (ami az indukciós feltevés miatt lehetséges).

Hasonlóan járhatunk el akkor is, amikor w első betűje b , csak ilyenkor a levezetés az $S \Rightarrow \mathbf{bSaS}$ szabállyal indul. \square

Nézzünk egy kicsit bonyolultabb nyelvtant!

4. Példa.

$$R \rightarrow XRX \mid S \quad (1-2)$$

$$S \rightarrow \mathbf{aTb} \mid \mathbf{bTa} \quad (3-4)$$

$$T \rightarrow XTX \mid X \mid \varepsilon \quad (5-7)$$

$$X \rightarrow \mathbf{a} \mid \mathbf{b} \quad (8-9)$$

Ez is környezetfüggetlen nyelvtan, ahol a kezdő változó az R . Az alábbi levezetésben az aláhúzott rész jelöli a következőként alkalmazott szabály bal oldalát, a nyíl feletti szám a szabály sorszámát.

$$\begin{aligned} R &\stackrel{1}{\Rightarrow} \underline{XRX} \stackrel{1}{\Rightarrow} \underline{XXRX} \stackrel{2}{\Rightarrow} \underline{XXSXX} \stackrel{8}{\Rightarrow} \underline{\mathbf{aXSXX}} \stackrel{9}{\Rightarrow} \underline{\mathbf{abSXX}} \\ &\stackrel{8}{\Rightarrow} \underline{\mathbf{abSaX}} \stackrel{8}{\Rightarrow} \underline{\mathbf{abSaa}} \stackrel{3}{\Rightarrow} \underline{\mathbf{abaTbaa}} \stackrel{5}{\Rightarrow} \underline{\mathbf{abaXTbaa}} \stackrel{7}{\Rightarrow} \underline{\mathbf{abaXXbaa}} \\ &\stackrel{9}{\Rightarrow} \underline{\mathbf{ababXbaa}} \stackrel{9}{\Rightarrow} \underline{\mathbf{ababbaa}} \end{aligned}$$

Tehát a kapott szó $\mathbf{ababbaa} \in L(G)$.

A fenti levezetés során több választásunk is volt, hogy melyik változót melyik szabály alapján helyettesítsük.

Egy levezetés sokszor jobban áttekinthető ha a lépéseket egy fába rendezzük.

4. Definíció. Legyen G egy környezetfüggetlen nyelvtan és x egy szó. Az x levezetési fája G -ben egy gyökeres fa, melyben

- a gyökér a kezdő változóval,
- minden nem levél csúcs egy-egy változóval,
- minden levél pedig Σ egy-egy elemével (vagy ε -nal) van címkézve.
- Ha egy A csúcs gyerekei balról jobbra olvasva B_1, B_2, \dots, B_k , akkor a nyelvtannak van $A \rightarrow B_1 B_2 \dots B_k$ szabálya. (Itt $B_i \in \Sigma \cup V \cup \{\varepsilon\}$.)
- A levelek balról jobbra olvasva éppen az x szót adják.

A definícióból világos, hogy egy $x \in L(G)$ szó tetszőleges levezetéséből lehet levezetési fát készíteni, és a levezetési fából is kiolvasható legalább egy levezetés.

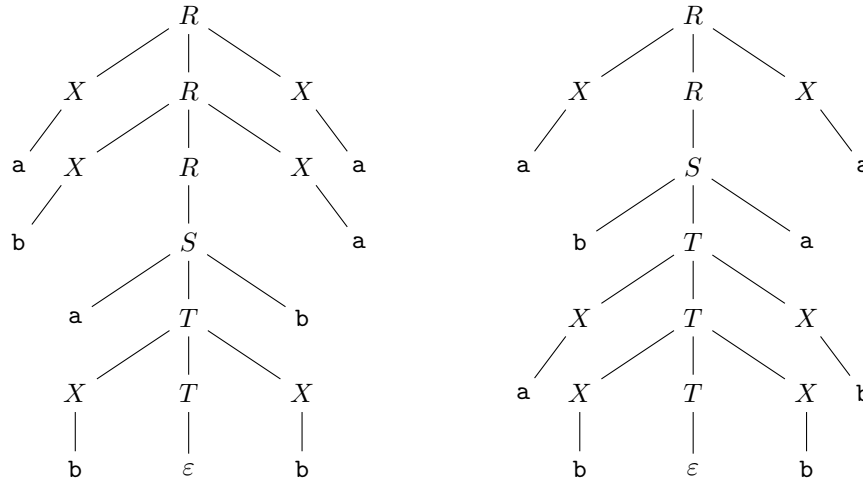
Fontos azonban megjegyezni, hogy míg a levezetés egyértelműen meghatározza a fát, visszafelé ez nem igaz, általában egy levezetési fából ugyanannak a szónak több levezetése is kiolvasható.

5. Definíció. Egy $x \in L(G)$ szó bal-levezetése egy olyan levezetés, amikor minden lépésben a γ_i elejéhez legközelebbi változót helyettesítjük egy megfelelő nyelvtani szabály alapján.

Erre már igaz, hogy egy levezetési fából egyetlen bal-levezetés olvasható ki.

5. Példa. Az előző példában leírt levezetéshez tartozó levezetési fa. Ebből többféle levezetés is leolvasható, a bal-levezetés szabályai sorrendben: 1, 8, 1, 9, 2, 3, 5, 9, 7, 9, 8, 8.

Jobb oldalt egy ugyanehhez a szóhoz tartozó másik levezetési fa látható.



6. Definíció. Egy $w \in \Sigma^*$ szó egyértelműen levezethető a G nyelvtanból, ha G -ben csak egy levezetési fája van.

A G nyelvtan egyértelmű, ha G -ből minden $w \in L(G)$ szó egyértelműen levezethető.

Az L nyelv egyértelmű, ha létezik egyértelmű nyelvtana.

Ezek szerint az előző példabeli szó nem egyértelműen levezethető, hiszen két különböző levezetési fája is van. Így persze a nyelvtan sem egyértelmű.

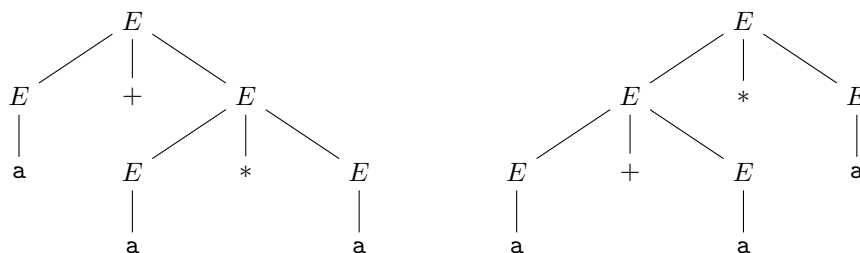
2. Megjegyzés. Az egyértelműen levezethetőség fenti definíciója ekvivalens azal, hogy a szó bal-levezetése egyértelmű.

Lássunk most egy fontos példát, az aritmetikai kifejezések nyelvét. Az egyszerűség kedvéért csak összeadást és szorzást fogunk benne használni, de egészíthető további műveletekkel is.

$$E \rightarrow E + E \mid E * E \mid (E) \mid a \quad (1)$$

Itt E az egyetlen változó, az ábécé elemei pedig $+$, $*$, a , valamint a nyitó és csukó zárójel.

Ez egy nem egyértelmű nyelvtan, hiszen például az $a + a * a$ kifejezéshez két különböző levezetési fa is tartozik,



3. Megjegyzés. Ha erre a két fára nem mint levezetési fákra, hanem mint a kifejezés kiértékelésének módját megadó fákra gondolunk, akkor látszik, hogy míg az első megfelel az aritmetikai kifejezések szokásos kiértékelésének (előbb a szorzást végezzük el, utána az összeadást) a második „rossz sorrendben” végzi a műveleteket.

1. Tétel. Az aritmetikai kifejezésekre adott fenti G egyszerű (1) nyelvtan nem egyértelmű, de az általa generált $L(G)$ nyelv egyértelmű nyelv.

Bizonyítás vázlat: Az előbb már láttuk, hogy a nyelvtan nem egyértelmű. A nyelv egyértelműségéhez mutatnunk kell egy G' egyértelmű nyelvtant, amire $L(G') = L(G)$. Legyen G' a következő:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid a \end{aligned}$$

Világos, hogy a G' nyelvtannal levezethető aritmetikai kifejezések levezethetők az eredeti nyelvtanból is.

Azt kell megmutatni, hogy ha $w \in L(G)$, akkor $w \in L(G')$ is teljesül, sőt a G' -beli levezetési fája egyértelmű.

Ezt a w hossza szerinti indukcióval mutathatjuk meg. Ha $|w| = 1$, akkor csak $w = a$ lehet, és ez egyedül az $E \Rightarrow T \Rightarrow F \Rightarrow a$ lépésekkel kapható meg a G' nyelvtanban.

Hosszabb szavakra azt kell észrevenni, hogy ha vannak zárójelen kívüli $+$ jelek, akkor először ezeket kell generálni (sorrendben visszafelé) az első szabály segítségével, utána a zárójelen kívüli $*$ jeleket, majd a zárójelekben levő kifejezéseket.

□

4. Megjegyzés. Vegyük észre, hogy ebben a módosított nyelvtanban ha a levezetési fát kiértékelési fának tekintjük, akkor a műveletek sorrendje is a szokásos lesz.