

A GÉPI /MONDAT/ANALIZIS NÉHÁNY KÉRDÉSÉRŐL

Ábrahám Sámuel

Napjaink egyik legfontosabb tudományága az információ feldolgozás és viszsza nyeres elmélete és gyakorlata. Modern tudományunk és technikánk mind nagyobb és nagyobb mennyiségű információ feldolgozását és visszanyerését igényli és pedig meghatározott, legtöbbször rövid időn belül. Ezeknek az igényeknek a kielégítésére az emberi erő egyedül már nem elégséges. Ezért égető sürgősséggel merül fel az információ feldolgozás és visszanyerés gépesítésének problémája, vagyis különböző gépek felhasználása az információ feldolgozás és visszanyerés folyamatában. A folyamatra engedélyezett idő határozza meg, hogy mechanikus vagy elektronikus gépek használata a legcélszerűbb, és egyáltalán nem kell arra törekedni, hogy minden egyes esetben függetlenül az adott körülményektől, mindent a leggyorsabb elektronikus gépekre bizzunk. Mindamellet, hogy elméletileg az egész folyamat egységes gépesítése tűnik optimálisnak, gyakorlatilag egy megfelelő ember-gép rendszer kidolgozása /ez leginkább talán a gépi fordítás esetében látszik/ a legkedvezőbb megoldás.

Mivelhogy a mai társadalomban az információ legnagyobb része természetes nyelvek segítségével rögzítődik, az információ feldolgozás és visszanyerés /egyik/ legfontosabb része a természetes nyelveken írt /ritkán: kiejtett/ szövegek feldolgozása. Ezért a jelen tanulmányban csak a természetes nyelveken írt szövegek gépi feldolgozásával foglalkozunk, pontosabban az ilyen szövegekből kapott információ gépi visszanyerésének egyes kérdéseivel. A kapcsolatban levő mondatok / = szövegek/ gépi feldolgozásának kérdései csak a legutóbbi időben kerültek az érdeklődés középpontjába, mert az egyszerűbb kérdés, a mondatok gépi feldolgozása sem nyert még elégséges megoldást.

A mondatok gépi feldolgozása elméletileg egyenlő a mondatok formális feldolgozásának fogalmával. A formális feldolgozást már pontosan meg lehet határozni: egy feldolgozás formális, ha csak a mondatot alkotó jelek formáját és mondaton belüli sorrendjét veszi tekintetbe /ezt és csak ezt tudja a gép érzékelni/. Hangsúlyoznunk kell, hogy formális feldolgozás révén /mindamellet, hogy határai mindinkább tágulnak/ lehetetlen a mondatban tartalmazott teljes információt visszanyerni. /A következőkben a "mondat formális feldolgozása a benne tartalmazandó információ visszanyerésének céljából"

kifejezés helyett a "mondat formális analízise" kifejezést használjuk, és amikor ez nem vezethet félreértéshez, a "formális" szót elhagyjuk./

Az utóbbi idők kutatásai arra utalnak, hogy a formális analízis könnyebben megvalósítható, ha először egy olyan /formális/ rendszert építünk fel, amely előállítja /generálja/ a mondatokat, és utólag "fordítjuk" meg ezt a rendszert. Ezért észlelhetjük napjainkban a generatív rendszerek /generatív grammatikák/ iránti nagy érdeklődést.

Ezeknek a grammatikáknak a szerkesztését és tanulmányozását N. CHOMSKY /1/ indította el, /és folytatja számottevő eredménnyel./ CHOMSKY szerint egy g e n e r a t í v m o d e l l /most már világos, hogy a "grammatika"szó használata nem teljesen indokolt/ egy olyan formális utasítás-rendszer /device/, amely az adott nyelv összes mondatait /és csak azokat/ állítja elő, mégpedig úgy, hogy a generálás mikéntje alapján már egyértelműen meghatározhatjuk a /generált/ mondat összes formális tulajdonságait. Ezt úgy is mondhatjuk, hogy a generálás folyamatában minden mondathoz ez az utasítás-rendszer adekvát analízist kapcsoljon. Természetesen a fenti elgondolás csak akkor válik ténylegesen meghatározássá, ha pontosan megfogalmazzuk, mit értünk "formális utasítás-rendszer"-en, és melyek a mondat "összes" formális tulajdonságai. Az első megtehetjük. Formális utasítás-rendszer alatt egy $G = /V, \Sigma, F/$ hármast értünk, ahol V egy véges szótár, amely véges számú szimbólumot tartalmaz. Σ egy részhal-maza V-nek, F pedig egy véges halmaz, amely / φ, ψ / kettéseket tartalmaz, ahol φ, ψ a V-ben tartalmazott szimbólumok sorozatai. Az F-be tartozó kettéseket törvényeknek nevezzük. Egy ilyen utasítás-rendszer a következőképpen generál:

Azt mondjuk, hogy σ_i -re alkalmazható a $f_k / = / \varphi_k, \psi_k //$ törvény, ha $/a/ \sigma_i$ olyan szimbólumok sorozata, amelyek V elemei, $/b/ \sigma_i$ tartalmazza φ_k és $/c/$ teljesül egy adott C_k feltétel. Ha $/a/, /b/, /c/$ teljesül, akkor f_k /egyszeri/ alkalmazása σ_i -re abból áll, hogy leírjuk a σ_j szimbólum sorozatot, amely σ_i -től csak annyiban különbözik, hogy σ_i -ban /az első/ φ_k helyett ψ_k -t írunk. Ebben az esetben azt mondjuk, hogy σ_i /közvetlenül/ generálja σ_j -t. Terminális derivációnak egy olyan $\sigma_1, \dots, \sigma_n$ sorozatot nevezünk, amelyben $/a/ \sigma_1$ generálja σ_{i+1} -et /minden $1 \leq i \leq n-1$ -re/, $/b/ \sigma_i$ eleme Σ -nak és $/c/ \sigma_n$ nem generál semmit. σ_n -t mondatnak, és a mondatok halmazát, LG-t, a G utasítás-rendszer által generált nyelvnek nevezzük.

Egyszerű példája a magyar nyelv kis részét generáló modellnek a következő:

$$G_m : V = \{ S, NP, VP, A, N, az, a, ember, könyvet, felveszi \}$$

$$= \{ S \}$$

$$F = /1/ \quad /S, NP VP/ \quad /4/ \quad /A, \begin{Bmatrix} az \\ a \end{Bmatrix} /C4/$$

$$/2/ \quad /NP, A N/ \quad /5/ \quad /N, \begin{Bmatrix} ember \\ könyvet \end{Bmatrix} /C5/$$

$$/3/ \quad /VP, V NP/ \quad /6/ \quad /V, felveszi/.$$

Mielőtt valamit is generálnánk ezzel a modellel, két megjegyzést kell tennünk. A {} zárójel azt jelenti, hogy az egyik vagy a másik szimbólumot választjuk ki, a C₄ /vagy C₅/-ben megfogalmazott feltételek szerint. A feltételeket /C₄, C₅/ célszerű a törvények leírásánál hozzájuk kapcsolni. /Az /1/, /2/, /3/, /6/ törvények alkalmazása nincs semmilyen feltételhez kötve./

Ha C₄ és C₅-öt megfelelően fogalmazzuk meg, akkor G_m a következőképpen generál egy magyar mondatot /két szimbólum-sorozat közé írjuk az alkalmazott törvény számát/.

S /1/, NP VP /2/, A N VP /3/, A N V NP //,

A N V A N /4/, Az N V A N /5/, az ember

V A N /6/, az ember felveszi A N /5/,

az ember felveszi a könyvet /4/, a z e m b e r f e l v e s z i a k ö n y v e t .

Sokkal nehezebb a feladatunk, ha az "összes formális tulajdonságok" fogalmát akarjuk pontosan meghatározni. Általában ezt a fogalmat csak magyarázzák, az adott nyelvet beszélő egyének "nyelvi intuíciójával", /amely fogalom ugyancsak nincs meghatározva/.

Az F-ben tartalmazott törvények és a hozzájuk kapcsolt C_k feltételek formájának különféle meghatározásával különböző generatív modelleket kapunk. Közülük érdemlegesek az úgynevezett

/1/ Kontextus mentes mondat szerkezet modellje,

/2/ Kontextust tekintetbe vevő mondat szerkezet modellje,

/3/ Transzformációs modell. /2/

Meg kell jegyeznünk, hogy míg az /1/ és /2/ modellek esetében sikerült pontosan meghatározni a törvények és feltételek formáját, a transzformációs modell esetében ezt még kielégítő módon elérni nem sikerült.

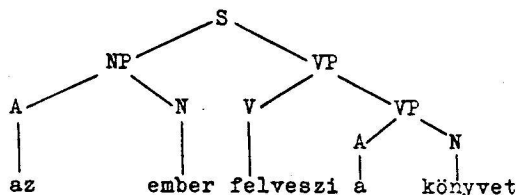
A jelen tanulmány értelmében részünkre különösen fontos az, hogy a modell ne csak generálja az adott természetes nyelvet, /vagyis generatív ereje legyen/, hanem hogy a generálás folyamata minden egyes mondatához egy megfelelő analízist is kapcsoljon /vagyis explikatív ereje is legyen/.

Az /1/ és /2/-es modellekben minden generált mondatához egy címkével ellátott ágrajzot /labelled-tree, P-marker/ kapcsolunk úgy, hogy a terminális deriváció tagjait sorrendben egymás alá írjuk és vonallal kötjük össze azokat a szimbólumokat /két egymásután következő sorban/, amelyek együtt egy f_k törvényt alkotnak.

A fenti példa esetében,

a z e m b e r f e l v e s z i a k ö n y v e t

- mondatnak a következő P-marker felel meg



CHOMSKY-nak sikerült kimutatni, hogy az /1/ és /2/ modellek /a fenti megfogalmazásokban/ még ha generálják is az adott természetes nyelvet /ez eddig nem sikerült ezt a kérdést teljesen tisztázni/ az általuk biztosított analízis nem adekvát. Ezért az egyetlen helyes modellnek a transzformációs modellt tekintik, amelynek egy része egy /2/ típusú modell. Ebben a modellben először véges számú un. m a g m o n d a t o t /pontosabban majdnem mondatot/ generálnak, amelyekből kiindulva generálják az összes többi mondatokat. Ebben a modellben egy mondat analízise, ha magmondatról van szó, egy címkével ellátott ágrajz, ha nem magmondatról van szó, akkor több címkével ellátott ágrajz /amelyek között vannak a generálásnál felhasznált magmondatok címkével ellátott ágrajzai/.

Ideje, hogy áttérjünk a fenti modellek "megfordítása" lehetőségének tanulmányozására. Az /1/ és /2/ modelleket meg lehet fogalmazni mint fél-asszociatív-rendszerek, és ezen az alapon minden egyes ilyen modellnek egyértelműen meg lehet feleltetni egy Markov-féle normál algoritmust. /3/ A Markov-féle normál algoritmusokat pedig könnyű "megfordítani": Minden egyes \mathcal{C} Markov-féle normál algoritmus esetében fel lehet építeni /bizonyos feltételek mellett egyértelműen/ egy olyan \mathcal{C} Markov-féle normál algoritmust, hogy ha $\mathcal{C}/P/ = R$, akkor $\mathcal{C}/R/ = P$. A probléma részletezését l.: /4/ és /5/-ben. Mivel a Markov-féle normál algoritmusok gépileg megvalósíthatók, a fenti modellek esetében a gépi analízis teljes egészében elvégezhető.

Amint már említettem, az /1/ és /2/ típusú modellek, CHOMSKY megfogalmazásában, nem megfelelőek a természetes nyelvek leírására, a transzformációs modell, /amelyről feltételezhetjük, hogy megfelelő/ pedig sajnos még nincs formális szempontból kielégítően megfogalmazva, így "megfordításuk" kérdését még csak fel sem lehet vetni. Ebből a helyzetből három kiút lehetséges.

CHOMSKY és követői azon fáradoznak, hogy a transzformációs modell egy formális szempontból kielégítő meghatározását adják. Ezek fontos, és szerintünk eredménnyel kecsegtető munkálatok, de eddig végleges eredmények még nem születtek.

A kutatók egy kisebb csoportja újabb fajta modelleket próbál felépíteni, amelyek eredményesebbek lesznek mind az /1/ és /2/-es típusnál. /Az eddigi eredmények viszonya a transzformációs grammatikához nincs tisztázva./ Korai lenne még eredményekről beszélni.

Szerintünk egy harmadik út is járható. Tul erősnek tűnik

CHOMSKY-nak az a követelménye, hogy egy generatív modell úgy legyen felépítve, hogy már a generálás alapján adja ki a megfelelő analízist. Ha ezt a két követelményt, a generálás és a megfelelő analízis felépítésének követelményét, szétválasztjuk, akkor az előttünk álló feladat megoldhatóvá válik. Eddig sikerült bebizonyítani, hogy ha adva van egy transzformációs grammatika /az eddigi megfogalmazások értelmében/, akkor fel lehet építeni egy vele generatív ekvivalens, vagyis egy pontosan ugyanazt a nyelvet generáló, /2/-es típusu grammatikát.

Ebben a grammatikában újabb formális műveletek bevezetésével minden egyes, már generált mondathoz egy olyan analízist lehet kapcsolni, amely egyenértékű a transzformációs grammatikában kapott analízissel. Ezt az új modellt //2/-es típusu grammatika + új analízis rendszer/ már meg lehet fogalmazni mint Markov-féle normál algoritmust, és ebből következően meg lehet "fordítani".

Az ilyen újfajta, általunk generatív-analitikus-nak nevezett modellek elméletén és gyakorlati kérdésein KIEFER Ferenc tudományos munkatárssal együtt dolgozunk a Magyar Tudományos Akadémia Számítástechnikai Központja keretében, és az eredményekről a közeljövőben részletesen beszámolunk.



BIBLIOGRÁFIA

- /1/ CHOMSKY, N.: Syntactic Structures. The Hague, 1963. Mouton and Co. 118 p.
- /2/ KIEFER F.: Matematikai nyelvészet. Bp. 1964. OMKDK 179 p.
- /3/ MARKOV, A.A.: Teoria algoritmov. Trudi. mat. inst. AN SSR, XLII, Moszkva, 1964. AN SSR 375 p.
- /4/ ÁBRAHÁM, S.: A formal study of generative grammars 1. Computational linguistics /Budapest/ 1964, II.sz. 5-21.p.
- /5/ ÁBRAHÁM, S.: A formal study of generative grammars 2. Computational linguistics /Budapest/ 1965. IV.sz. /nyomás alatt/

oOo

Ш. Абрахам

Некоторые вопросы машинного анализа /предложений/

В статье рассматриваются известные до сих пор методы формального анализа /графически изображенных/ предложений в рамках порождающих грамматик. После указания недостатков этих методов, показывается, что они вытекают из /слишком сильного/ требования, чтобы анализ был полностью и исключительно обоснован на процессе порождения предложения. Если от этого требования отказаться, то порождающую грамматику можно дополнить системой определений таким образом, чтобы удовлетворительный формальный анализ стал полностью осуществим. Такая работа проводится в рамках Вычислительного центра ВАН. Более подробно новый метод будет изложен в будущей статье.

ÁBRAHÁM, S.: Some questions of the mechanical analysis of sentences

In the present paper methods of formal language analysis /within generative grammars/ are presented. After discussing their shortcomings, it is shown that latter are due to the /to strong/ restriction that the analysis should be based solely on the way the sentence has been generated. If we do not respect this restriction, generative grammar can be completed by a conceptual machinery so that adequate formal analysis becomes possible even on the base of lower-type generative grammars. Studies on such systems are carried out at the Computing Centre of the Hungarian Academy of Sciences, and will be reported in a future paper in more detail.

=

= a