

Θεωρία Υπολογισμού

Γλώσσες και Γραμματικές Χωρίς Συμφραζόμενα

2003-04

Γλώσσες και Γραμματικές Χωρίς Συμφραζόμενα

- Γλώσσες και γραμματικές χωρίς συμφραζόμενα.
- Αυτόματα στοίβας.
- Ιδιότητες των γλωσσών χωρίς συμφραζόμενα.
- Εφαρμογές στη συντακτική ανάλυση.

Γραμματικές Χωρίς Συμφραζόμενα

- Οι γραμματικές χωρίς συμφραζόμενα μελετήθηκαν αρχικά από γλωσσολόγους (για παράδειγμα τον Noam Chomsky).
- Παράδειγμα:

$Sentence \rightarrow NounPhrase Verb$
 $NounPhrase \rightarrow Adjective NounPhrase$
 $NounPhrase \rightarrow Noun$
 $Noun \rightarrow boy$
 $Adjective \rightarrow little$
 $Verb \rightarrow run$

Γραμματικές Χωρίς Συμφραζόμενα

- Ο BNF (Backus-Naur-Form) συμβολισμός που χρησιμοποιούμε για να περιγράψουμε το συντακτικό γλωσσών προγραμματισμού έχει ουσιαστικά τις ίδιες έννοιες με τις γραμματικές χωρίς συμφραζόμενα.
- Παράδειγμα:

$Expression ::= Expression + Expression$
 $Expression ::= Expression * Expression$
 $Expression ::= (Expression)$
 $Expression ::= id$

- Ο ορισμός των παραπάνω αριθμητικών εκφράσεων δεν μπορεί να γίνει με κάποια κανονική έκφραση.

Ο Ρόλος των Γραμματικών

- Η χρήση γραμματικών προσφέρει σημαντικά πλεονεκτήματα στους σχεδιαστές γλωσσών προγραμματισμού, τους υλοποιητές μεταγλωττιστών και τους προγραμματιστές:
 - Μια γραμματική δίνει **ακριβή ορισμό** του συντακτικού μιας γλώσσας.
 - Κάποιες κλάσεις γραμματικών επιτρέπουν την αυτόματη κατασκευή αποδοτικών συντακτικών αναλυτών που αναγνωρίζουν αν ένα πρόγραμμα είναι συντακτικά σωστό (και αν δεν είναι υποδεικνύουν το λάθος που υπάρχει).
 - Η εισαγωγή μιας καινούριας έννοιας σε μια γλώσσα προγραμματισμού (και η υλοποίησή της) είναι πιο εύκολη αν υπάρχει ήδη μια γραμματική για τη γλώσσα (και οι μεταγλωττιστές είναι βασισμένοι στη γραμματική αυτή).

Γραμματικές Χωρίς Συμφραζόμενα

- **Ορισμός. Γραμματική χωρίς συμφραζόμενα** (context-free grammar) είναι μια τετράδα $G = (V, \Sigma, R, S)$ όπου
 - V είναι ένα αλφάβητο
 - Σ , το **σύνολο των τερματικών**, είναι ένα υποσύνολο του V . Τα στοιχεία του $V \setminus \Sigma$ ονομάζονται **μη τερματικά**.
 - R , το **σύνολο των κανόνων (rules or productions)**, είναι ένα πεπερασμένο υποσύνολο του $(V \setminus \Sigma) \times V^*$
 - S , το **αρχικό σύμβολο**, είναι ένα στοιχείο του $V \setminus \Sigma$.
- Ένας κανόνας $(A, u) \in R$ γράφεται συνήθως σαν $A \rightarrow_G u$ ή $A \rightarrow u$.

Παράδειγμα

$$G = (V, \Sigma, R, S)$$

- όπου
 - $V = \{S, a, b\}$
 - $\Sigma = \{a, b\}$
 - $R = \{S \rightarrow aSb, S \rightarrow e\}$
- Οι γραμματικές χωρίς συμφραζόμενα είναι **παραγωγοί γλωσσών**.
- Η παραπάνω γραμματική παράγει τη (μη κανονική) γλώσσα $\{a^n b^n : n \geq 0\}$.

Ορισμοί

- Η σχέση \Rightarrow_G (παράγει σε ένα βήμα):
Αν $u, v \in V^*$ τότε $u \Rightarrow_G v$ (ή $u \Rightarrow v$) ανν υπάρχουν $x, y, v' \in V^*$ και $A \in V \setminus \Sigma$ έτσι ώστε $u = xAy$, $v = xv'y$ και $A \rightarrow v'$.
- Η σχέση \Rightarrow_G^* (παράγει) είναι η ανακλαστική και μεταβατική κλειστότητα της \Rightarrow_G .
- Κάθε ακολουθία της μορφής
$$w_0 \Rightarrow_G w_1 \Rightarrow_G \dots \Rightarrow_G w_n$$
λέγεται **παραγωγή** (derivation) της w_n από την w_0 στην G .
- Μια τέτοια παραγωγή λέμε ότι έχει **μήκος** n (ή αποτελείται από n βήματα).

Ορισμοί

- Αν $G = (V, \Sigma, R, S)$ είναι μια γραμματική χωρίς συμφραζόμενα τότε η γλώσσα που παράγεται από την G είναι

$$L(G) = \{w \in \Sigma^* : S \Rightarrow_G^* w\}$$

- Μια γλώσσα L ονομάζεται **γλώσσα χωρίς συμφραζόμενα** (context-free language) αν $L = L(G)$ όπου G είναι μια γραμματική χωρίς συμφραζόμενα.
- Δυο γραμματικές G_1 και G_2 λέγονται **ισοδύναμες** αν $L(G_1) = L(G_2)$.

Παράδειγμα

$$G = (V, \Sigma, R, E)$$

- όπου

$$V = \{E, T, F, \text{id}, +, *, (,)\}$$

$$\Sigma = \{\text{id}, +, *, (,)\}$$

και

$$R = \{ E \rightarrow T + E, \quad E \rightarrow T,$$

$$T \rightarrow F * T, \quad T \rightarrow F,$$

$$F \rightarrow (E), \quad F \rightarrow \text{id} \}$$

- Ποιά είναι η γλώσσα που παράγεται από την παραπάνω γραμματική;

Γλώσσες Προγραμματισμού και Γραμματικές

- Όταν χρησιμοποιούμε γραμματικές χωρίς συμφραζόμενα για να ορίσουμε το συντακτικό γλωσσών προγραμματισμού τότε:
 - Τα τερματικά σύμβολα παριστάνουν τα tokens της γλώσσας προγραμματισμού (δηλαδή τις καλά ορισμένες συμβολοσειρές π.χ. `if`, `x12`, `==` στην C).
 - Τα μη τερματικά σύμβολα παριστάνουν τα υπόλοιπα δομικά στοιχεία των γλωσσών προγραμματισμού.

Παράδειγμα

$$G = (V, \Sigma, R, S)$$

- όπου

$$V = \{S, (,)\}$$

$$\Sigma = \{ (,) \}$$

και

$$R = \{ S \rightarrow e, \quad S \rightarrow SS, \quad S \rightarrow (S) \}$$

- Ποιά είναι η γλώσσα που παράγεται από την παραπάνω γραμματική;

Παράδειγμα

- Δύο παραγωγές της G είναι

$$S \Rightarrow SS \Rightarrow S(S) \Rightarrow S((S)) \Rightarrow S(() \Rightarrow (S)() \Rightarrow ()()$$

και

$$S \Rightarrow SS \Rightarrow (S)S \Rightarrow ()S \Rightarrow ()(S) \Rightarrow ()((S)) \Rightarrow ()()$$

- Έτσι βλέπουμε ότι μπορεί η ίδια συμβολοσειρά να έχει περισσότερες από μία παραγωγές σε μία γραμματική χωρίς συμφραζόμενα.

Αριστερότερες Παραγωγές

- Κάθε παραγωγή της μορφής

$$w_0 \Rightarrow_G w_1 \Rightarrow_G \dots \Rightarrow_G w_n$$

λέγεται **αριστερότερη παραγωγή** (leftmost derivation) αν το μη τελικό σύμβολο που αντικαθίσταται σε κάθε βήμα είναι το αριστερότερο.

- **Παράδειγμα:** $G = (V, \Sigma, R, S)$ όπου

$$V = \{S, (,)\}$$

$$\Sigma = \{ (,) \}$$

και

$$R = \{ S \rightarrow e, \quad S \rightarrow SS, \quad S \rightarrow (S) \}$$

Παράδειγμα

- Η παραγωγή

$$S \Rightarrow SS \Rightarrow (S)S \Rightarrow ()S \Rightarrow ()(S) \Rightarrow ()((S)) \Rightarrow ()()$$

είναι αριστερότερη, αλλά η

$$S \Rightarrow S\underline{S} \Rightarrow S(\underline{S}) \Rightarrow S((S)) \Rightarrow S(() \Rightarrow (S)() \Rightarrow ()()$$

δεν είναι.

Αριστερότερες Παραγωγές

- Γράφουμε $\alpha \Rightarrow_G^L \beta$ όπου $\alpha, \beta \in V^*$ αν $\alpha = \alpha_1 A \alpha_2$, $\beta = \alpha_1 \gamma \alpha_2$, $A \rightarrow \gamma$ είναι ένας κανόνας της G και $\alpha_1 \in \Sigma^*$.
- Η ανακλαστική και μεταβατική κλειστότητα της \Rightarrow_G^L είναι η \Rightarrow_G^{*L} .

Αριστερότερες Παραγωγές

- **Θεώρημα:** Για κάθε γραμματική χωρίς συμφραζόμενα $G = (V, \Sigma, R, S)$ και κάθε συμβολοσειρά $w \in \Sigma^*$, $S \Rightarrow_G^* w$ αν και μόνο αν $S \Rightarrow_G^{*L} w$.

Απόδειξη:

- (Αν) Προφανές.
- (Μόνο αν) Έστω

$$w_0 \Rightarrow_G w_1 \Rightarrow_G \dots \Rightarrow_G w_n$$

μια παραγωγή όχι υποχρεωτικά αριστερότερη, όπου $w_0 = S$ και $w_n = w \in \Sigma^*$.

Η παραπάνω παραγωγή μπορεί να μετασχηματιστεί σε μια αριστερότερη παραγωγή.

Συνέχεια της Απόδειξης

- Έστω k ο μικρότερος αριθμός τέτοιος ώστε το βήμα $w_k \Rightarrow_G w_{k+1}$ δεν είναι η αριστερότερη παραγωγή του w_{k+1} από το w_k .
- Τότε η w_k θα είναι της μορφής $\alpha A \beta B \gamma$ όπου $\alpha \in \Sigma^*$, $\beta, \gamma \in V^*$, $A, B \in V \setminus \Sigma$, η w_{k+1} θα είναι της μορφής $\alpha A \beta \delta \gamma$ και υπάρχει ο κανόνας $B \rightarrow_G \delta$.
- Εφόσον $w_n \in \Sigma^*$, ισχύει $k+1 < n$ και υπάρχει κάποιο επόμενο βήμα όπου αντικαθίσταται και το A . Έτσι έστω l ο μικρότερος αριθμός, μεγαλύτερος του k έτσι ώστε $w_l = \alpha A \epsilon$ για κάποιο $\epsilon \in V^*$ και $w_{l+1} = \alpha \zeta \epsilon$, όπου $A \rightarrow_G \zeta$. Τότε $\beta B \gamma \Rightarrow_G^* \epsilon$ σε $l-k$ βήματα.

Συνέχεια της Απόδειξης

- Αρα μπορούμε να εναλλάξουμε τις χρήσεις των κανόνων $B \rightarrow_G \delta$ και $A \rightarrow_G \zeta$ ως εξής.

$$\begin{array}{ll} w_0 \Rightarrow_G^{*L} w_k = \alpha A \beta B \gamma & k \text{ βήματα} \\ \Rightarrow_G^L \alpha \zeta \beta B \gamma & 1 \text{ βήμα} \\ \Rightarrow_G^L \alpha \zeta \epsilon & l-k \text{ βήματα} \\ \Rightarrow_G^L w_n & n-l-1 \text{ βήματα} \end{array}$$

- Η νέα παραγωγή έχει μήκος n και αρχίζει με τουλάχιστο $k+1$ βήματα αριστερής παραγωγής.

Παράδειγμα

- Η παραγωγή

$$S \Rightarrow SS \Rightarrow S(S) \Rightarrow S((S)) \Rightarrow S(() \Rightarrow (S)() \Rightarrow ()()$$

μπορεί να γίνει αριστερότερη εφαρμόζοντας το παραπάνω θεώρημα διαδοχικά ως εξής:

– Βήμα 1:

$$S \Rightarrow SS \Rightarrow (S)S \Rightarrow (S)(S) \Rightarrow (S)((S)) \Rightarrow (S)() \Rightarrow ()()$$

– Βήμα 2:

$$S \Rightarrow SS \Rightarrow (S)S \Rightarrow ()S \Rightarrow ()(S) \Rightarrow ()((S)) \Rightarrow ()()$$

Συντακτικά Δέντρα

- Τα συντακτικά δέντρα (parse trees) είναι μια γραφική μέθοδος αναπαράστασης παραγωγών.
- Θεωρήστε τη γραμματική $G = (V, \Sigma, R, S)$ όπου

$$V = \{S, (,)\}$$

$$\Sigma = \{(),\}$$

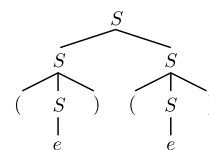
$$R = \{ S \rightarrow e, S \rightarrow SS, S \rightarrow (S) \}$$

Συντακτικά Δέντρα

- Η αριστερότερη παραγωγή

$$S \Rightarrow SS \Rightarrow (S)S \Rightarrow (S)(S) \Rightarrow ()(S) \Rightarrow ()()$$

μπορεί να παρασταθεί από το ακόλουθο συντακτικό δέντρο:

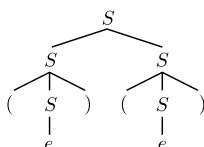


Συντακτικά Δέντρα

- Αλλά και η παραγωγή

$$S \Rightarrow SS \Rightarrow S(S) \Rightarrow (S)(S) \Rightarrow (S)() \Rightarrow ()()$$

μπορεί να παρασταθεί από το ίδιο συντακτικό δέντρο:



Συντακτικά Δέντρα

- **Ορισμός.** Ένα δέντρο λέγεται **συντακτικό δέντρο** για την γραμματική χωρίς συμφραζόμενα $G = (V, \Sigma, R, S)$ αν:
 - Οι κόμβοι του δέντρου συμβολίζονται από στοιχεία του συνόλου $V \cup \{e\}$.
 - Η ρίζα του δέντρου συμβολίζεται από το αρχικό σύμβολο S .
 - Κάθε εσωτερικός κόμβος συμβολίζεται με ένα μη τερματικό σύμβολο $A \in V \setminus \Sigma$.
 - Αν ο κόμβος A έχει παιδιά Q_1, \dots, Q_n διαβασμένα από δεξιά προς αριστερά, τότε

$$A \rightarrow Q_1 \dots Q_n$$

είναι ένας κανόνας του R .

- Κάθε κόμβος που συμβολίζεται με e είναι φύλλο και είναι το μοναδικό παιδί του πατέρα του.

- Το **προϊόν** (yield) ενός συντακτικού δέντρου είναι η συμβολοσειρά που προκύπτει με παράθεση από τα αριστερά προς τα δεξιά των συμβόλων της γραμματικής που βρίσκονται στα φύλλα του δέντρου.

Διττές Γραμματικές

- **Ορισμός:** Μια γραμματική G λέγεται **διττή** αν έχει παραπάνω από μια διαφορετικές αριστερότερες παραγωγές για την ίδια συμβολοσειρά $w \in L(G)$.

- **Παράδειγμα:** Έστω η γραμματική $G = (V, \Sigma, R, E)$ όπου

$$V = \{E, T, F, \text{id}, +, *\}, \quad \Sigma = \{\text{id}, +, *\}$$

και

$$R = \{ E \rightarrow E + E, E \rightarrow E * E, E \rightarrow \text{id} \}$$

Διττές Γραμματικές

- Η γραμματική αυτή έχει δύο διαφορετικές αριστερότερες παραγωγές για τη συμβολοσειρά $\text{id} + \text{id} * \text{id}$:

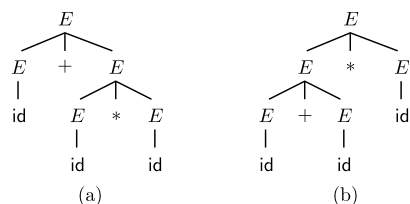
$$E \Rightarrow E + E \Rightarrow \text{id} + E \Rightarrow \text{id} + E * E \Rightarrow \text{id} + \text{id} * E \Rightarrow \text{id} + \text{id} * \text{id}$$

και

$$E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow \text{id} + E * E \Rightarrow \text{id} + \text{id} * E \Rightarrow \text{id} + \text{id} * \text{id}$$

Διττές Γραμματικές

- Άρα έχουμε και δύο διαφορετικά συντακτικά δέντρα:



Σχέση με τις Κανονικές Γλώσσες

- Το σύνολο των κανονικών γλωσσών είναι υποσύνολο του συνόλου των γλωσσών χωρίς συμφραζόμενα.

- **Παράδειγμα:** Η γραμματική $G = (\{S, a\}, \{a\}, R, S)$ όπου

$$R = \{ S \rightarrow aS, S \rightarrow e \}$$

παράγει τη γλώσσα a^* .

Κανονικές Γραμματικές

- **Ορισμός.** Μια γραμματική χωρίς συμφραζόμενα $G = (V, \Sigma, R, S)$ λέγεται **κανονική** αν

$$R \subseteq (V \setminus \Sigma) \times \Sigma^*((V \setminus \Sigma) \cup \{e\})$$

- **Διαισθητικά:** Μια γραμματική χωρίς συμφραζόμενα λέγεται **κανονική** αν κάθε κανόνας είναι της μορφής

$$N_1 \rightarrow wN_2$$

ή

$$N_1 \rightarrow w$$

όπου w είναι μια συμβολοσειρά αποτελούμενη από τερματικά σύμβολα και N_1, N_2 είναι μη τερματικά σύμβολα.

Παράδειγμα

- Η γραμματική

$$G = (V, \Sigma, R, S)$$

όπου

$$V = \{S, A, B, a, b\}$$

$$\Sigma = \{a, b\}$$

$$R = \{ S \rightarrow bA, S \rightarrow aB, S \rightarrow e,$$

$$A \rightarrow abaS, B \rightarrow babS \}$$

είναι κανονική.

- Η γλώσσα που παράγεται από την γραμματική αυτή είναι $L(G) = (baba \cup abab)^*$.

Κανονικές Γραμματικές

- **Θεώρημα.** Μια γλώσσα είναι κανονική αν μπορεί να παραχθεί από μια κανονική γραμματική.

Απόδειξη: (Μόνο αν)

- Έστω μια κανονική γλώσσα που γίνεται δεκτή από το αυτόματο $M = (K, \Sigma, \delta, s, F)$. Τότε $L(M) = L(G)$ όπου $G = (V, \Sigma, R, S)$ είναι η κανονική γραμματική που ορίζεται ως εξής:

$$- V = \Sigma \cup K$$

$$- S = s$$

$$- R = \{ q \rightarrow ap : \delta(q, a) = p \} \cup \{ q \rightarrow e : q \in F \}$$

Συνέχεια της Απόδειξης

- Οι κανόνες της G έχουν σχεδιαστεί ακριβώς για να μιμούνται τις κινήσεις του M .

- Δηλαδή για κάθε $\sigma_1, \dots, \sigma_n \in \Sigma$ και $p_0, \dots, p_n \in K$,

$$(p_0, \sigma_1 \dots \sigma_n) \vdash_M (p_1, \sigma_2 \dots \sigma_n) \vdash_M \dots \vdash_M (p_n, e)$$

ανν

$$p_0 \Rightarrow_G \sigma_1 p_1 \Rightarrow_G \sigma_1 \sigma_2 p_2 \dots \Rightarrow_G \sigma_1 \dots \sigma_n p_n$$

26

28

30

32

25

27

29

31

Συνέχεια της Απόδειξης

- Ας υποθέσουμε ότι $w \in L(M)$. Τότε $(s, w) \vdash_M^* (p, e)$ για κάποιο $p \in F$. Τότε $s \Rightarrow^* wp$ και εφόσον $p \rightarrow e$ είναι επίσης κανόνας της G , θα ισχύει $s \Rightarrow^* w$ και $w \in L(G)$. Άρα $L(M) \subseteq L(G)$.
- Ας υποθέσουμε ότι $w \in L(G)$. Τότε $s \Rightarrow_G^* w$. Ο κανόνας που χρησιμοποιήθηκε στο τελευταίο βήμα της παραγωγής πρέπει να ήταν της μορφής $p \rightarrow e$, για κάποιο $p \in F$, έτσι ώστε $s \Rightarrow_G^* wp \Rightarrow_G^* w$. Όμως τότε $(s, w) \vdash_M^* (p, e)$ και $w \in L(M)$. Άρα $L(G) \subseteq L(M)$.
- Τελικά $L(G) = L(M)$.

Συνέχεια της Απόδειξης

(Αν)

- Έστω $G = (V, \Sigma, R, E)$ μία κανονική γραμματική. Ορίζουμε ένα μη ντετερμινιστικό πεπερασμένο αυτόματο M έτσι ώστε $L(M) = L(G)$.
- Έστω $M = (K, \Sigma, \Delta, s, F)$, όπου

$$K = (V \setminus \Sigma) \cup \{f\} \quad (f \text{ είναι ένα νέο σύμβολο})$$

$$s = S$$

$$F = \{f\}$$

$$\Delta = \{(A, w, B) : A \rightarrow wB \in R \text{ και } A, B \in V \setminus \Sigma \text{ και } w \in \Sigma^*\}$$

$$\cup \{(A, w, f) : A \rightarrow w \in R \text{ και } A \in V \setminus \Sigma \text{ και } w \in \Sigma^*\}.$$

Συνέχεια της Απόδειξης

- Οι κινήσεις του M μιμούνται τις παραγωγές της G . Για κάθε $A_1, \dots, A_n \in V \setminus \Sigma, w_1, \dots, w_n \in \Sigma^*$,

$$A_1 \Rightarrow_G w_1 A_2 \Rightarrow_G w_1 w_2 A_3 \dots$$

$$\dots \Rightarrow_G w_1 w_2 \dots w_{n-1} A_n \Rightarrow_G w_1 w_2 \dots w_{n-1} w_n$$

ανν

$$(A_1, w_1 \dots w_n) \vdash_M (A_2, w_2 \dots w_n) \vdash_M \dots \vdash_M (A_n, w_n) \vdash_M (f, e).$$

Συνέχεια της Απόδειξης

- Επομένως, αν $w \in L(G)$ τέτοια ώστε $w \in \Sigma^*$ και $S \Rightarrow_G^* w$, τότε $(S, w) \vdash_M^* (f, e)$ και έτσι $w \in L(M)$. Άρα $L(G) \subseteq L(M)$.
- Αντίστροφα, αν $w \in L(M)$ τέτοια ώστε $(S, w) \vdash_M^* (f, e)$ τότε $S \Rightarrow_G^* w$, έτσι ώστε $w \in L(G)$. Άρα $L(M) \subseteq L(G)$.
- Τελικά $L(M) = L(G)$.

Μελέτη

- Κεφάλαια 3.1-3.2 από το βιβλίο:
Harry R. Lewis και Χρίστος Χ. Παπαδημητρίου, Στοιχεία Θεωρίας Υπολογισμού, Τεχνικό Επιμελητήριο Ελλάδας, 1992.