

# 13068. Teoria d'Autòmats i Llenguatges Formals

## 1. Introducció

Francesc J. Ferri  
Dept. d'Informàtica. Universitat de València

4 d'octubre de 2002

TALF (1. Introducció)

## Definicions preliminars

**símbol (etiqueta, caràcter):** unitat lògica indivisible.

**alfabet (vocabulari):** conjunt finit de símbols  $x$

$\Sigma_1 = \{a, b, c\}$ ,  $\Sigma_2 = \{0, 1, 2, \dots, 9\}$ ,  $\Sigma_3 = \{\text{vertader, fals, indefinit}\}$

**cadena (sentència, paraula, frase):** successió finita de símbols d'un determinat alfabet.

$x_1 = aaabbbba$

$x_2 = 1024$

$x_3 = \text{vertader} \cdot \text{vertader} \cdot \text{indefinit}$

## [Definicions preliminars]

**longitud (talla):** nombre de símbols d'una cadena

$$|x_1| = 7, |x_2| = 4, |x_3| = 3$$

Notació:  $|x|_A = \{\text{nombre de símbols de } A \subseteq \Sigma \text{ en } x\}$

$$|x_1|_a = |x_1|_{\{a\}} = 4, |x_1|_c = 0, |x_1|_{\{b,c\}} = 3$$

**concatenació:** juxtaposició de dues cadenes. Operació interna dins del conjunt de totes les possibles paraules sobre un determinat alfabet.

$$x \cdot y = xy, \quad aaa \cdot ab = aaaab$$

**cadena buida (nul·la,  $\varepsilon$ ):** cadena que no conté cap símbol. Element neutre de la concatenació.

## [Definicions preliminars]

**prefix:**  $y$  és prefix de  $x$  si  $\exists z$  tal que  $x = yz$ .

**sufix:**  $y$  és sufix de  $x$  si  $\exists z$  tal que  $x = zy$ .

**factor (subcadena, infix):**  $y$  és factor de  $x$  si  $\exists z, t$  tals que  $x = zyt$ .  
Els prefixs i sufixs en són casos particulars.

**factor (prefix, sufix) propi de  $x$ :** si és diferent de  $\varepsilon$  i de  $x$ .

**factorització de  $x$ :** descomposició d'una cadena en un nombre finit de factors,  $x = x_1 \cdot x_2 \cdots x_k$ .

## [Definicions preliminars]

**llenguatge:** qualsevol conjunt de cadenes.

**talla (grandària, cardinalitat) d'un llenguatge:** nombre de cadenes que conté.

**Llenguatges “especials”**  $\{\varepsilon\}$ , i  $\emptyset$

**Exemples de llenguatges**  $\{a, aa, aaa, \dots\}$ ,  $\{0, 10, 20\}$  i  $\{\langle \text{fals} \rangle \langle \text{fals} \rangle, \varepsilon\}$ , les talles dels quals són  $\infty$ , 3 i 2, respectivament

## Cadenes sobre un alfabet

$P(\Sigma)$  és el conjunt format per *totes* les cadenes sobre un alfabet.

$$P(\{0, 1\}) = \{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, \dots\}$$

Aquesta ordenació de les cadenes s'anomena **lexicogràfica**.

El conjunt  $P(\Sigma)$  té estructura de *monoide* respecte de la concatenació (p. associativa, elem. neutre).

$P(\Sigma)$ : monoide lliure

$P(\Sigma) - \{\varepsilon\}$ : semigrup lliure

## Operacions sobre cadenes

**Potència:**  $x^i = \underbrace{x \cdots x}_i$  .  
 $i$  vegades

- $x^0 = \varepsilon$
- $x^i x^j = x^{i+j}$
- $x^{i+1} = x^i \cdot x = x \cdot x^i$
- $|x^i| = i \cdot |x|$

**Inversió (reflexió):**  $x^{-1} = a_n \cdots a_2 a_1$  si  $x = a_1 \cdots a_n$  i  $a_i \in \Sigma$

$(\text{pepito})^{-1} = \text{otipep}$

## Operacions sobre llenguatges

S'hi poden aplicar qualssevol operacions definides sobre conjunts

$$L_1 \cup L_2 \quad L_1 \cap L_2 \quad \overline{L_1} \quad L_1 - L_2$$

**concatenació:**  $L_1 \cdot L_2 = \{z \mid z = x \cdot y, x \in L_1, y \in L_2\}$

$\mathcal{L}(\Sigma)$ : conjunt de tots els llenguatges sobre  $\Sigma$  (també  $2^{P(\Sigma)}$ ), té estructura de monoide respecte de  $\cdot$  i  $\cup$ .

Els elements neutres són  $\{\varepsilon\}$  i  $\emptyset$ , resp.

Se l'anomena binoide lliure

## Exemples (concatenació)

- $L \cdot \emptyset = \emptyset$
- $L \cdot \{\varepsilon\} = L$
- $\{a^n | n \geq 0\} \cdot \{b^m | m \geq 0\} = \{a^n b^m | n \geq 0, m \geq 0\}$
- $\{a^n | n \geq 0\} \cdot \{a^n | n \geq 0\} = \{a^n | n \geq 0\}$

$$L_1 = \{a, b\}^* \quad L_2 = a \cdot \{a, b\}^*$$

$$L_1 \subset L_2?$$

$$|L_1| < |L_2|?$$

## [Operacions sobre llenguatges]

**Potència**  $L^i = \underbrace{L \cdot L \cdots L}_{i \text{ vegades}}$

$$L^i = \{x_1 \cdot x_2 \cdots x_i | x_k \in L, k = 1, \dots, i\}$$

$$L^0 = \{\varepsilon\}$$

$$L^{i+1} = LL^i = L^iL$$

$$L^{i+j} = L^iL^j = L^jL^i$$

**Inversió (reflexió):**  $L^{-1} = \{x^{-1} | x \in L\}$

## [Operacions sobre llenguatges]

**Tancament (clausura) de Kleene**  $L^* = \bigcup_{i=0}^{\infty} L^i$

$$\Sigma^* = P(\Sigma)$$

**Tancament positiu:**  $L^+ = \bigcup_{i=1}^{\infty} L^i$

$$\begin{aligned} L^* &= L^+ \cup \{\varepsilon\} \\ L^+ &= L^* - \{\varepsilon\} \text{ només si } \varepsilon \notin L! \\ L^+ &= LL^* = L^*L \end{aligned}$$

**Quocient:**  $L_1/L_2 = \{x \mid \exists y \in L_2 : xy \in L_1\}$

Prefixs de cadenes de  $L_1$  el corresponent sufix del qual està en  $L_2$ .

## Exemples (tancaments)

- $L^+L^+ = L^2L^* = LL^*L = L^*L^2$
- $\emptyset^* = \{\varepsilon\}$
- $\emptyset^+ = \emptyset$
- $\{a^n \mid n \geq 0\}^* = \{a^n \mid n \geq 0\}$
- $\{a^n b^m \mid n \geq 0, m \geq 0\}^* = \{a, b\}^*$

## Exemples (quocients)

Siga  $L_a = \{a^n | n \geq 0\}$ ,  $L_b = \{b^n | n \geq 0\}$ ,  $L_{ab} = \{a^n b^m | n \geq 0, m \geq 0\}$

- $L/\emptyset = \emptyset$
- $L/\{\varepsilon\} = L$
- $L_a/L_b = L_a/\{\varepsilon\} = L_a$
- $L_{ab}/L_b = L_{ab}$
- $L_{ab}/\{ab\} = L_a$
- $L_{ab}/L_a = L_{ab}$

## [Operacions sobre llenguatges]

**Substitució:** correspondència que associa un llenguatge a un símbol.

$$f : \Sigma \longrightarrow 2^{\Gamma^*}$$

$$\text{Podem definir } \begin{cases} f'(\varepsilon) = \varepsilon \\ f'(ax) = f(a)f'(x) \end{cases}$$

$$f' : \Sigma^* \longrightarrow 2^{\Gamma^*}$$

$$\text{També podem definir } f''(L) = \{f'(x) | x \in L\}$$

$$f'' : 2^{\Sigma^*} \longrightarrow 2^{\Gamma^*}$$

No distingirem entre  $f$ ,  $f'$  i  $f''$

## Exemple (substitució)

Siguen  $C$  i  $I$  dos llenguatges sobre l'alfabet  $\Sigma = \{A - z, 0 - 9, +\}$ .

$$C = \{1 - 9\} \cdot \{0 - 9\}^* \quad I = \{A - z\} \cdot \{A - z, 0 - 9\}^*$$

Siga  $L$  un llenguatge sobre l'alfabet  $\Gamma = \{c, i, +\}$ .

$$L = (\{c\} \cup \{i\}) \cdot (\{+c\} \cup \{+i\})^*$$

podem definir  $f(c) = C$ ,  $f(i) = I$ ,  $f(+) = +$ .

Aleshores  $f(L)$  representa sumes arbitràries de identificadors i constants enteres.

## [Operacions sobre llenguatges]

**Homomorfisme:** cas particular en què el llenguatge associat a un símbol consta d'una única cadena.

$$\begin{aligned} h &: \Sigma \longrightarrow \Gamma^* \\ h' &: \Sigma^* \longrightarrow \Gamma^* \\ h'' &: 2^{\Sigma^*} \longrightarrow 2^{\Gamma^*} \end{aligned}$$

**Homomorfisme invers:**

$$\begin{aligned} h^{-1} &: \Gamma^* \longrightarrow 2^{\Sigma^*} \quad !! \\ & \quad h^{-1}(x) = \{w \in \Sigma^* \mid h(w) = x\} \end{aligned}$$

De la mateixa manera:

$$h^{-1}(L) = \{w \in \Sigma^* \mid h(w) \in L\}$$



## Exemples (homomorfismes)

Siga  $h(a) = 0$  i  $h(b) = 00$

$$\begin{aligned} h^{-1}(x) &= \{w \in \Sigma^* \mid h(w) = x\} \\ h^{-1}(L) &= \{w \in \Sigma^* \mid h(w) \in L\} \end{aligned}$$

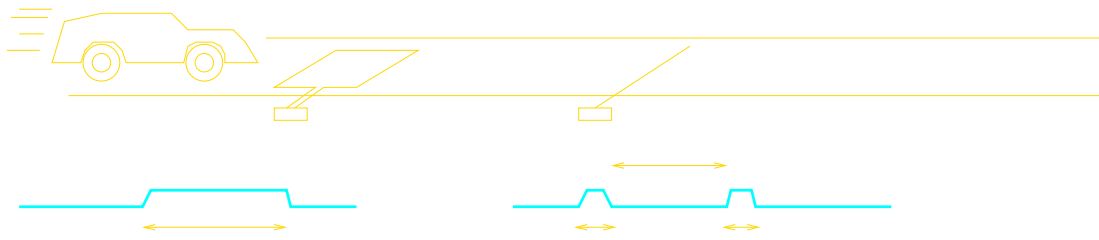
- $h(ab) = 000$
- $h^{-1}(000) = \{ab, ba, aaa\}$
- $h^{-1}(1) = \emptyset$
- $h^{-1}(\varepsilon) = \{\varepsilon\}$
  
- $h(L_{ab}) = \{0^n \mid n \geq 0\}$
- $h(L_b) = \{0^{2n} \mid n \geq 0\}$
- $h^{-1}(\{0^{2n} \mid n \geq 0\}) = \text{exercici!}$

## Exemples de llenguatges

- $L_1 = \{a^n \mid n \geq 0\}$
- $L_2 = \{a^n b^m \mid n \geq 0, m \geq 0\}$
- $L_3 = \{w \in \{a, b\}^* : |w|_a = |w|_b\}$
- $L_4 = \{w \in \{a, b\}^* : |w| \leq 2\} = \{\varepsilon, a, b, aa, ab, ba, bb\}$
- $L_5 = \{w \in \{a, b\}^* : w = w^{-1}\}$  (palíndromes)
- $L_6 = \{w \in \{a, b\}^* : w = x \cdot x^{-1}\}$  (palíndromes parelles)

EXERCICIS:  $L_4^*$ ,  $L_2 \cap L_3$ ,  $L_3/L_5$ ,  $L_5/L_6$ ,  $L_5/L_1$ , ...

## Més exercicis



Siga  $20 \leq T \leq 25$ ,  $2 \leq R \leq 3$ ,  $15 \leq E \leq 20$ .

Descriure formalment els llenguatges  $L_1$  i  $L_2$  que representen cadenes de 0 i 1 corresponent a que ha passat un cotxe donades les mesures anteriors.

Descriure'ls de la forma més compacta possible fent servir els mecanismes que s'han vist fins ara.

## i més

Volem detectar pàgines web (cadenes sobre l'alfabet ASCII) que continguin paraules amb el prefix propi **sex-**.

Exemple: sextet, sexagesimal, (Atenció: la pàgina sencera és una cadena!)

Descriure formalment de forma compacta el llenguatge corresponent.

## i més encara

Descriure formalment el llenguatge sobre l'alfabet  $\{a, (, )\}$  format per cadenes ben parentitzades com

$$(((a)(a))(a))a$$

- $L = \{w : |w|_{(} = |w|_{)}\}$ ?
- $L = \{w : |w|_{(} = |w|_{)} \wedge |x|_{(} \geq |x|_{)}, \forall x : w = xy\}$ ?
- $L = \{w : w = (x), x \in L \vee w = xy, x \in L, y \in L\}$ ?

## Inducció i recursió

L'última definició no té sentit però dóna la idea d'una definició inductiva per a  $L$

1.  $a \in L$ .
2.  $xy \in L$  si  $x \in L, y \in L$ .
3.  $(x) \in L$  si  $x \in L$ .

Aquesta definició es pot llegir recursivament:

Les cadenes de  $L$  són: o 1) una  $a$ , o 2) la concatenació de dues cadenes (més petites) de  $L$ , o 3) una cadena (més petita) de  $L$  envoltada per parèntesis.

## Gramàtiques

La definició recursiva anterior mijantçant regles:

1.  $\langle \text{Cadena} \rangle \rightarrow a$
2.  $\langle \text{Cadena} \rangle \rightarrow \langle \text{Cadena} \rangle \langle \text{Cadena} \rangle$
3.  $\langle \text{Cadena} \rangle \rightarrow (\langle \text{Cadena} \rangle)$

El símbol  $\langle \text{Cadena} \rangle$  representa qualsevol cadena de  $L$ . És una variable. Les regles expressen com pot estar formada aquesta cadena en funció d'altres cadenes.

Es pot dir que  $\langle \text{Cadena} \rangle$  es reescriu com a  $(\langle \text{Cadena} \rangle)$  (r. 3)

que es reescriu com a  $(\langle \text{Cadena} \rangle \langle \text{Cadena} \rangle)$  (r. 2)

que es reescriu com a  $(aa)$  (r. 1 dues vegades).

## Gramàtiques. definició

$$G = \langle V_T, V_N, S, P \rangle$$

- $V_T = \Sigma$  és l'alfabet o vocabulari terminal.
- $V_N$  és l'alfabet no terminal.
- $P$  és un conjunt de **produccions**  
 $\alpha \rightarrow \beta : \beta \in V^* \text{ i } \alpha \in V^* V_N^+ V^* \text{ (} V = V_T \cup V_N \text{)}.$
- $S \in V_N$  **axioma**.

## [Gramàtiques. definició]

**derivació directa**  $x \xrightarrow[G]{y}$   $\exists u \rightarrow v \in P \wedge x = \alpha u \beta, y = \alpha v \beta.$

**derivació** (tancament transitiu de  $\xrightarrow[G]{}$ )  $x \xrightarrow[G]{*} y$

$$x = y \vee x \xrightarrow[G]{a_1} \xrightarrow[G]{a_2} \cdots \xrightarrow[G]{a_n} y$$

**formes sentencials**  $x \in V^* : S \xrightarrow[G]{*} x$

**sentències** formes sentencials de  $V_T$

**llenguatge generat per  $G$**   $L(G) = \{w \in V_T^* \mid S \xrightarrow[G]{+} w\}$

## Exemples. Gramàtiques

- $S \rightarrow SS|(S)|a$   
 $S \Rightarrow S\underline{S} \Rightarrow S\underline{S}S \Rightarrow S(\underline{S})S \Rightarrow S(S\underline{S})S \Rightarrow S(S(S))S \xrightarrow{*} a(a(a))a$
- $S \rightarrow S + S|S * S|(S)|x$   
 $\underline{S} \Rightarrow \underline{S} * S \Rightarrow (\underline{S}) * S \Rightarrow (\underline{S} + \underline{S}) * \underline{S} \xrightarrow{+} (x + x) * x$
- $S \rightarrow 0S|A|\varepsilon$   
 $A \rightarrow 1A|\varepsilon$   
 $S \Rightarrow 0S \Rightarrow 00S \Rightarrow 00A \Rightarrow 001A \Rightarrow 0011A \Rightarrow 00111A \Rightarrow 001111$

## La jerarquia de Chomsky

### tipus 3. Regulars

$$A \rightarrow aB \quad \text{o} \quad A \rightarrow a \qquad a \in V_T, A, B \in V_N$$

### tipus 2. Incontextuals

$$A \rightarrow \alpha \qquad A \in V_N, \alpha \in V^*$$

### tipus 1. Contextuals

$$xAy \rightarrow x\beta y \quad \text{o} \quad S \rightarrow \varepsilon \qquad x, y \in V^*, A \in V_N, \beta \in V^+$$

### tipus 0. No restringides

## Jerarquia de Chomsky. Exemples

$$\text{tipus 3: } S \rightarrow 0A|1A|\varepsilon \\ A \rightarrow 0S|1S$$

$$\text{tipus 2: } S \rightarrow 0S0|1S1|0|1|\varepsilon$$

$$\text{tipus 1: } S \rightarrow T|\varepsilon \\ T \rightarrow aTBD|abD \quad CD \rightarrow BD \\ DB \rightarrow CB \quad bB \rightarrow bb \\ CB \rightarrow CD \quad D \rightarrow c$$

$$\text{tipus 0: } S \rightarrow aSAc|abc|\varepsilon \\ cA \rightarrow Ac \\ bA \rightarrow bb$$

## Gramàtica de tipus 0?

$$\begin{array}{l|l} S \rightarrow aSAc|abc|\varepsilon & S \rightarrow aSAc|abc|\varepsilon \\ cA \rightarrow Ac & cA \rightarrow cX \quad cX \rightarrow AX \\ bA \rightarrow bb & AX \rightarrow Ac \quad bA \rightarrow bb \end{array}$$

$$\begin{aligned} S &\Rightarrow aSAc \Rightarrow aaSAcAc \xrightarrow{*} a^{n-1}S(Ac)^{n-1} \Rightarrow \\ &\Rightarrow a^nb c(Ac)^{n-1} = a^nb(cA)^{n-1}c \xrightarrow{*} \boxed{a^nb(Ac)^{n-1}c} \\ &= a^nbAc(Ac)^{n-2}c \Rightarrow a^nb b c(Ac)^{n-2}c = \\ a^nb b(cA)^{n-2}cc &\xrightarrow{*} \boxed{a^nb b(Ac)^{n-2}cc} \\ \dots & \\ \xrightarrow{*} a^nb^{n-1}(cA)^1c^{n-1} &\xrightarrow{*} a^nb^{n-1}(Ac)c^{n-1} \Rightarrow a^nb^n c^n \end{aligned}$$

## Definicions alternatives

$$\alpha \rightarrow \beta : |\alpha| \geq |\beta| \implies \text{Contextual!}$$

$$xAy \rightarrow x\beta y \quad x, y \in V^*, A \in V_N, \boxed{\beta \in V^*} \implies \text{tipus 0!}$$

$$\begin{array}{l} A \rightarrow \alpha B \\ A \rightarrow \alpha \end{array} \quad A \in V_N, \alpha \in V_T^* \text{ ho considerarem regular.}$$

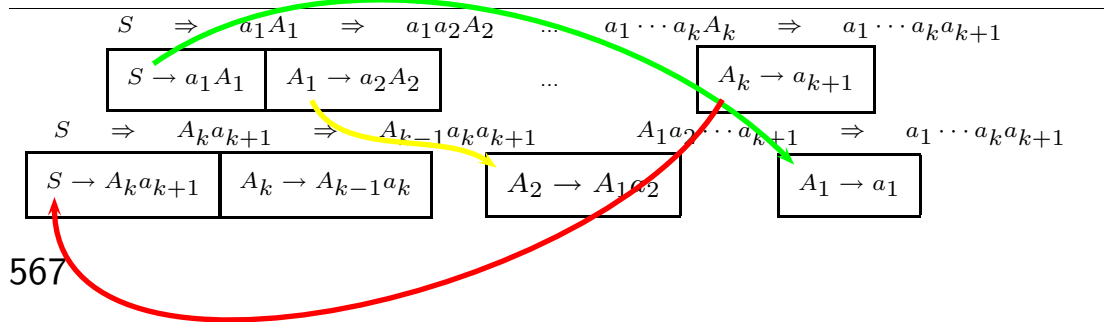
$$\begin{array}{l} A \rightarrow B\alpha \\ A \rightarrow \alpha \end{array} \quad A \in V_N, \alpha \in V_T^* \text{ regular per la dreta!}$$

$S \rightarrow 0S|S1|\varepsilon$  no és ni regular per la dreta ni per l'esquerra !!!

## Exemple. Gramàtiques regulars

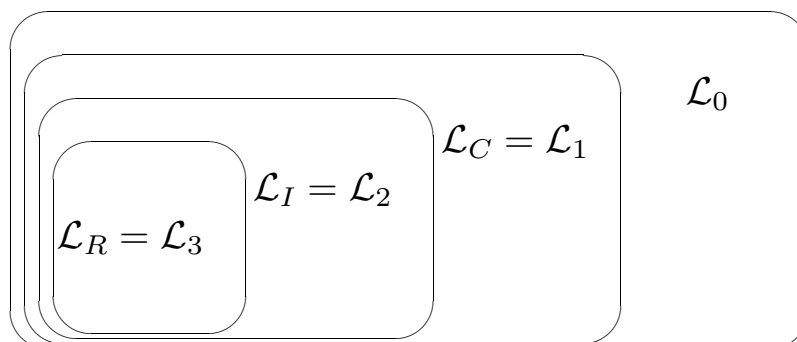
- $S \rightarrow S1|A|\varepsilon$  (per la dreta)  
 $A \rightarrow A0|\varepsilon$   
 $S \Rightarrow S1 \Rightarrow S11 \Rightarrow S111 \Rightarrow A111 \Rightarrow A0111 \Rightarrow A00111 \Rightarrow 00111$

- $S \rightarrow 0S|A|\varepsilon$  (per l'esquerra)  
 $A \rightarrow 1A|\varepsilon$   
 $S \Rightarrow 0S \Rightarrow 00S \Rightarrow 00A \Rightarrow 001A \Rightarrow 0011A \Rightarrow 00111A \Rightarrow 00111$



## Jerarquia de Chomsky. Llenguatges

Un llenguatge és de tipus  $N$  si existeix alguna gramàtica de tipus  $N$  que el genere.





## Exemples

$$\{a^n b^n c^n \mid n \geq 0\} \in \mathcal{L}_C \subseteq \mathcal{L}_0$$

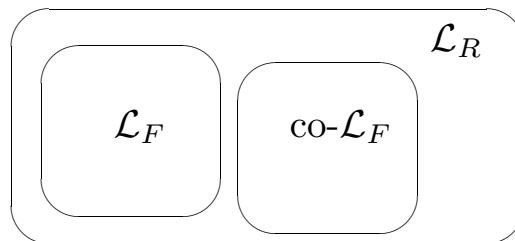
$$\{0^n 1^n \mid n \geq 0\} \in \mathcal{L}_I \subseteq \mathcal{L}_C \subseteq \mathcal{L}_0$$

**llenguatges finits:**  $|L| < \infty$

$\{a\}$

**llenguatges cofinits:**  $|\bar{L}| < \infty$

$\{a, b\}^+$



## Grans preguntes sense resposta ...

( ...de moment)

- $\mathcal{L}_R \subseteq \mathcal{L}_I \subseteq \mathcal{L}_C \subseteq \mathcal{L}_0$  però són pròpies aquestes inclusions?
- Hi ha llenguatges que no siguin  $\mathcal{L}_0$ ?
- On està la subclasse dels llenguatges coregulars?
- I els coincontextuals?
- Coincideix la noció intuïtiva de complexitat (estructural) d'un llenguatge amb els tipus de la jerarquia de Chomsky?
- Hi ha més classes dins de la jerarquia? I fora?