

## Capitolo C10

# stringhe e linguaggi formali

### Contenuti delle sezioni

- a. stringhe, linguaggi e monoidi liberi p. 2
- b. operazioni su stringhe e linguaggi p. 9
- c. morfismi tra linguaggi p. 15
- d. relazioni tra stringhe e derivazioni di linguaggi p. 19
- e. fattorizzazioni di stringhe p. 27
- f. parole di Lyndon p. 32

34 pagine

---

**C100.01** Questo è il primo dei capitoli del tomo C dedicato ai linguaggi formali, alle macchine sequenziali (o automi), alle considerazioni generali sulle attività computazionali e sulla loro complessità e alla introduzione della teoria dei codici e della crittografia.

Questi capitoli si sforzano di essere autonomi e il presente primo capitolo espone con una certa completezza le nozioni basilari su stringhe e linguaggi, già presentate come elementi iniziali dell'*esposizione* in B01 e riprese successivamente quando servivano a introdurre altre entità che si possono basare solo sulle stringhe e sulle liste [B06, B17].

Successivamente vengono presentati i linguaggi della gerarchia di Chomsky insieme alle macchine che li caratterizzano [C12-C16], fino a giungere alla macchine di Turing e alla computabilità [C20, C21].

Dopo aver introdotte alcune collezioni di linguaggi controllati con grammatiche e meccanismi meno classici [C23, C26], si esamina la fattorizzazione dei linguaggi [C30] e si espongono teorie più astratte sui linguaggi [C32-C35].

In C47 viene introdotto lo studio delle complessità computazionale e in C60 sono presentate alcune parole infinite.

Gli ultimi capitoli [C65, C66] sono dedicati alla introduzione delle probabilità discrete e dei codici.

## C10 a. stringhe, linguaggi e monoidi liberi

**C10a.01** In molte attività conoscitive e in molte manipolazioni materiali si devono esaminare sequenze finite di oggetti (più o meno composti) curandosi solo delle posizioni che essi occupano nelle sequenze stesse.

Gli oggetti componenti di queste sequenze vengono individuati da alcune caratteristiche che sono giudicate sostanziali in dipendenza dagli scopi per i quali gli oggetti stessi saranno utilizzati.

Sulla base di queste caratteristiche due oggetti di una sequenza possono essere giudicati diversi o indistinguibili; nel secondo caso si parla di repliche di un unico tipo di oggetto.

Tra le attività concernenti sequenze di oggetti vi sono le manipolazioni di elenchi di parole, di locuzioni, di frasi o di espressioni che fanno parte di lingue naturali o artificiali.

In questi casi gli oggetti componenti di sequenze, chiamati anche items di liste, vengono considerati elementari e si possono identificare con le loro scritte.

Per queste scritte in genere si auspica costituiscano dei validi significanti in grado di evitare dubbi identificativi e di facilitare la comprensione dei significati e dei ruoli.

Altre attività riguardano sequenze di oggetti sensibilmente complessi per i quali si chiede comunque che si possano rappresentare mediante scritte convenzionali che facciano riferimento a definizioni modellistiche anche articolate (che sul piano pratico possono essere soggette a discussioni).

Alcuni esempi: un elenco anagrafico, l'elenco dei processori costituenti un sistema multiprocessore, la sequenza dei residui di aminoacidi costituenti una proteina di medie dimensioni (10000 residui), un catalogo astronomico, un listino di prodotti tecnologici portatori di interrelazioni importanti per le applicazioni prevedibili.

In relazione agli oggetti ai quali si possono associare per le loro possibili applicazioni, per le sequenze che stiamo prospettando, per le loro componenti che sono considerate elementari e per gli oggetti tangibili dei quali queste componenti sono i modelli proponiamo, risp., i termini **sequenze rappresentative**, **items rappresentativi** e **oggetti rappresentati**.

**C10a.02** Accade che nelle varie aree disciplinari dei molti generi di sequenze rappresentative vengono effettuate e comunicate analisi nelle quali ciascun item, per quanto complesso possa essere l'oggetto tangibile che esso rappresenta, quando lo si incontra viene trattato solo in quanto segno (unità semi-ologica) che serve solo a distinguerlo dagli items da considerare diversi e ad assimilarlo agli items che si è convenuto di considerare indistinguibili.

Quando si dovesse andare più a fondo sugli oggetti tangibili rappresentati si dovrebbe ricorrere alle precisazioni delle relazioni fra oggetti rappresentati e items rappresentativi secondo il modello applicativo adottato.

Per completezza va detto che nelle scritte della tradizione pittografica e ideografica (in particolare nelle scritte della tradizione cinese) le unità simboliche hanno anche il compito di richiamare, almeno per sommi capi, gli oggetti rappresentati.

Conviene segnalare esplicitamente che l'ordine secondo il quale gli items compaiono nelle sequenze può essere collegato a motivi molto diversi: disposizione temporale, collocazione spaziale di qualche genere, ordinamento secondo vari tipi di valutazione e altro.

Questa varietà di situazioni corrisponde alla amplissima utilizzazione delle sequenze nelle attività e negli studi.

**C10a.03** Nell’ambito delle considerazioni sulle sequenze rappresentative gli items loro componenti si possono qualificare come segni che sono stati dotati o che potrebbero essere dotati di significati univoci, ossia potrebbero qualificarsi come contrassegni.

In un complesso di attività che riguarda un corpo di sequenze rappresentative risulta utile raggruppare i segni degli items in insiemi che per realismo dobbiamo dichiarare finiti ma che possiamo pensare estendibili quando interessano eventi conoscitivi o costruttivi che si sviluppano nel tempo.

Questi raggruppamenti di segni sono realizzati o realizzabili mediante elenchi nonripetitivi che quando sono opportune precisazioni semantiche si possono riconoscere come entrate della prima colonna di una opportuna “tabella esplicativa”. Questi elenchi nonripetitivi li chiamiamo **repertori di contrassegni**; in alcuni contesti è conveniente chiamarli **alfabeti**.

Per taluni complessi di attività questi repertori riguardano elenchi di poche decine di segni elementari. Tipici esempi sono l’elenco delle 26 lettere minuscole dell’alfabeto inglese, l’elenco dei 62 caratteri alfanumerici e l’elenco dei caratteri ASCII visualizzabili; questi elenchi presentano il vantaggio di essere pienamente gestibili dai più comuni strumenti informatico-telematici.

In questi casi i segni possono essere chiamati “caratteri” o “lettere”. e i loro insiemi “alfabeti elementari”.

Oltre agli alfabeti elementari sono in uso alfabeti di segni più elaborati, il più importante di questi è l’alfabeto *Unicode*.

Le sequenze finite di caratteri saranno dette **stringhe**; in alcuni ambiti disciplinari (algebra, logica) viene usato il termine **parole**. Ricordiamo che accanto alle parole, che sono sequenze di caratteri di lunghezza finita, sono studiate anche le parole infinite, che sono successioni di caratteri [C60].

**C10a.04** Segnaliamo alcuni alfabeti elementari particolarmente utilizzati:

il singoletto contenente solo un particolare segno che potrebbe essere  $\{ \mid \}$  e chiamato “tacca”;

l’alfabeto dei bits  $\{0, 1\}$ ;

l’insieme delle cifre decimali;

l’insieme delle lettere utilizzate in una lingua naturale (che potrebbe comprendere o meno lettere dotate di segni diacritici e potrebbe distinguere o meno tra lettere minuscole e maiuscole);

l’insieme dei caratteri di un sistema standardizzato di codifica come ASCII ( $w_i$ ) o Unicode ( $w_i$ ).

Ricordiamo per questi ultimi che tutti i caratteri ASCII si codificano con 7 o con 8 bits, mentre in Unicode i più basilari richiedono 16 bits e altri 32 e altri multipli di 16.

Va segnalato che in certe argomentazioni i caratteri ASCII o Unicode sono trattati come elementari, mentre in altre sono da considerare come stringhe sull’alfabeto dei bits.

Quale punto di vista scegliere dipende dai fini dello sviluppo argomentativo che si sta effettuando e a queste scelte conviene lasciare libertà ed elasticità ai singoli esecutori di procedure o ai singoli redattori di documenti.

**C10a.05** Per ragioni di praticità espositiva introduciamo il termine **insieme dei caratteri per alfabeti** ed il corrispondente simbolo  $\mathbb{A}$ .

Questo consente di servirsi di scritture come  $A \subset_{\varphi} \mathbb{A}$  e  $\{a_1, \dots, a_r\} \subset \mathbb{A}$  per introdurre semplici notazioni per alfabeti (come  $A$ ) e per caratteri (come gli  $a_r$ ).

$\mathbb{A}$  va considerato insieme finito che potrebbe essere individuato da una ampia indagine sopra le attività SciTech, cioè sulle ricerche e sulle realizzazioni di carattere scientifico e tecnologico; esso inoltre va considerato un insieme in continua evoluzione.

Consideriamo dunque un alfabeto  $A = \{a_1, \dots, a_n\} \subset_{\varphi} \mathbb{F}$ . Per la stringa esemplare su tale alfabeto usiamo la scrittura

$$w = \langle a_{j_1}, a_{j_2}, \dots, a_{j_r} \rangle \in A^r,$$

ove  $r$  è un intero positivo e  $a_{j_1}, \dots, a_{j_r}$  sono caratteri di  $A$ ;  $r$  viene detto **lunghezza della stringa**  $w$  e questa quantità si può denotare con  $|w|$ ,  $w^{\perp}(w)$  o  $\text{len}(w)$ .

**C10a.06** Chiamiamo **giustapposizione** o **concatenazione** delle stringhe  $w = \langle a_{j_1}, a_{j_2}, \dots, a_{j_r} \rangle$  e  $x = \langle a_{h_1}, a_{h_2}, \dots, a_{h_s} \rangle$  la stringa

$$w_1 x := \langle a_{j_1}, a_{j_2}, \dots, a_{j_r}, a_{h_1}, a_{h_2}, \dots, a_{h_s} \rangle.$$

Evidentemente questa operazione binaria è associativa: se  $w$ ,  $x$  e  $y$  denotano tre stringhe, la stringa ottenuta giustapponendo prima  $w$  ed  $x$ , poi  $w_1 x$  con  $y$  si assume coincidere con la stringa ottenuta in un primo momento con la giustapposizione di  $x$  e  $y$  ed in un secondo momento con la giustapposizione di  $w$  con  $x_1 y$ .

Queste due costruzioni quando si vogliono distinguere è ragionevole denotarle, risp., con  $(w_1 x)_1 y$  e con  $w_1 (x_1 y)$ .

Riconosciuta l'associatività della giustapposizione si può evitare la scrittura di alcune coppie di parentesi: invece di ciascuna delle precedenti scritture si può scrivere più semplicemente  $w_1 x_1 y$ .

Spesso si adotta anche la semplificazione consistente nel trascurare il segno che ricorda la giustapposizione e si scrive concisamente  $w x y$ .

In seguito svolgeremo varie considerazioni su stringhe e linguaggi di tono algebrico e per alcuni aspetti vicine a quelle più tradizionali riguardanti i campi numerici; per analogia spesso l'operazione di giustapposizione verrà chiamata anche **prodotto** e conseguentemente invece del segno “ $_1$ ” verrà usato il segno “ $\cdot$ ”.

Conviene rilevare esplicitamente che la assunzione della associatività per la giustapposizione discende dalla opportunità di utilizzare sul piano formale le stringhe per rappresentare i messaggi che si scambiano sistematicamente gli agenti matematico-informatici, messaggi la cui associatività risulta chiaramente giustificata.

Questo consente di descrivere, con efficacia, le stringhe come oggetti ottenuti mediante accostamento materiale di entità tangibili semplici (per esempio come affiancamento di caratteri tipografici) che assumono i ruoli delle sue successive componenti.

**C10a.07** Ogni stringa, potendosi considerare la giustapposizione dei suoi caratteri, può essere presentat con una scrittura come

$$w = a_{j_1} \cdot a_{j_2} \cdot a_{j_3} \cdot \dots \cdot a_{j_{r-1}} \cdot a_{j_r}.$$

Spesso il segno della giustapposizione, come altri segni attribuibili a qualche genere di “operazione prodotto”, per brevità si riduce un piccolo spazio tra i due successivi operandi (caratteri componenti) o anche si fa scomparire.

Vengono quindi adottate scritture come

$$w = a_{j_1} a_{j_2} a_{j_3} \dots a_{j_{r-1}} a_{j_r} = a_{j_1} \cdot a_{j_2} \cdot a_{j_3} \dots a_{j_{r-1}} \cdot a_{j_r}.$$

Si può quindi descrivere la giustapposizione di due stringhe  $w$  ed  $x$  come il processo di “scrivere” la  $w$  seguita immediatamente dalla  $x$  sopra un dispositivo in grado di registrare in successione i caratteri della  $w$  e dopo di questi i caratteri della  $x$ .

Tale dispositivo di registrazione sequenziale dei caratteri viene convenzionalmente chiamato **nastro**. Esso si può descrivere come un supporto che presenta una sequenza di sezioni elementari che chiamiamo **caselle**, ciascuna delle quali in grado di registrare un carattere dell'alfabeto con cui si lavora.

I nastri vengono presentati disposti orizzontalmente e si conviene di denotare come loro prima casella la più a sinistra, in sintonia con le modalità di scrittura e lettura che procedono da sinistra a destra (modalità prevalenti nei paesi “occidentali”).

Come si è visto in B17 e come vedremo in C21, risulta utile servirsi di nastri che, oltre a poter essere registrati da sinistra a destra, possono essere percorsi in avanti e all'indietro e avere le caselle modificate.

Inoltre in molti sviluppi argomentativi un nastro, cioè un oggetto tangibile e quindi finito, quando le circostanze lo richiedono, si consente possa essere dotato di ulteriori caselle da aggiungere sulla destra ed anche sulla sinistra.

A questo punto è opportuno segnalare che finora abbiamo denotato i caratteri dell'alfabeto di lavoro servendoci della fonte teletype e quindi di notazioni come  $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{a}_j$  o  $\mathbf{b}_{h_i}$ .

Tuttavia per il seguito ci riserviamo invece di usare notazioni come  $a$ ,  $b$ ,  $c$ ,  $a_j$  o  $b_{h_i}$  in quanto più agevolmente inseribili in espressioni matematiche elaborate.

**C10a.08** Per l'insieme delle stringhe su un alfabeto  $A$  si utilizza la notazione

$$A^+ := \bigcup_{r=1}^{+\infty} A^r .$$

$A^+$  può leggersi “A cross” e munito della giustapposizione dal punto di vista dell'algebra costituisce un semigruppato chiamato **semigruppato libero sull'alfabeto  $A$**  [B41a].

Si osserva che la giustapposizione è commutativa sse l'alfabeto ha un solo carattere, mentre per un alfabeto di due o più caratteri non lo è; ad esempio  $\mathbf{ab}, \mathbf{a} \neq \mathbf{a}, \mathbf{ab}$ .

Conviene osservare esplicitamente che la noncommutatività della giustapposizione è una caratteristica che rende le stringhe strumenti in grado di rappresentare gli eventi nel mondo reale, in quanto entità osservabili nel loro svolgersi nel tempo e strettamente vincolate dalle relazioni “prima di” e “dopo di”.

I semigruppato liberi della forma  $\langle A^+, \cdot \rangle$  non hanno unità; tuttavia a ciascuno di essi, come a ogni semigruppato [B41a], si può aggiungere una nuova entità che denotiamo con  $\mu$ , che considereremo stringa su ciascuno degli alfabeto utilizzabili e che chiamiamo **stringa muta** o **stringa vuota**, chiedendo che giustapposta a sinistra o a destra a qualsiasi stringa la lasci invariata:

$$w \cdot \mu = \mu \cdot w = w .$$

Essa quindi ha il ruolo di unità bilaterale per la giustapposizione.

Alla stringa muta si attribuisce la lunghezza 0,  $\mu^{-1} := 0$ ; inoltre per ogni alfabeto  $A$  si scrive  $A^0 := \{ \mu \}$ .

Si introduce allora la notazione

$$A^* := \bigcup_{r \in \mathbb{N}} A^r ,$$

da leggersi “A star” e si osserva che  $\langle A^*, \cdot, \mu \rangle$  costituisce un monoide del quale  $\mu$  è l'unità; esso è detto **monoide libero su un alfabeto  $A$** .

**C10a.09** Una stringa  $w \in A^r$  si può considerare una funzione del genere  $\lceil r \rceil \mapsto A$  oppure del genere  $\lceil r \rceil \mapsto A$  a seconda che si attribuiscono alle sue componenti, risp., indici come  $1, \dots, r$  o gli indici  $0, \dots, r-1$ .

Secondo le due convenzioni, si evidenziano i successivi caratteri componenti con le scritture, risp.,  $w = w(1)...w(r)$  oppure  $w = w(0)...w(r-1)$ ; nel seguito seguiremo prevalentemente la prima convenzione.

Quando risulterà opportuno scriveremo  $w =_A a_{j_1} a_{j_2} \dots a_{j_r}$  per segnalare sia l'uguaglianza delle stringhe a sinistra e a destra del segno "=", sia il fatto che ciascuno dei simboli  $a_{j_1}, a_{j_2}, \dots, a_{j_r}$  rappresenta un puro carattere dell'alfabeto  $A$  con cui si lavora.

Per lunghezza di una stringa  $w \in A^*$  si è definito l'intero naturale che esprime il numero delle sue componenti, cioè l'entità fornita dalla stringa ottenuta trasformando ogni carattere della  $w$  in in unico segno (per esempio  $\mid$ ). Si constata che

$$\forall w, x \in A^* : (w \mid x)^{\mid} = w^{\mid} + x^{\mid} .$$

In termini algebrici questa relazione dice che la funzione lunghezza costituisce un morfismo di ogni monoide libero  $\langle A^*, \mid, \mu \rangle$  sul monoide  $\langle \mathbb{N}, +, 0 \rangle$ .

Evidentemente si tratta di un epimorfismo che è biiettivo sse  $|A| = 1$ .

Le classi di equivalenza di questo morfismo sono per i diversi  $r \in \mathbb{N}$  gli insiemi  $A^r$  delle stringhe aventi la stessa lunghezza  $r$ .

Per ogni  $w \in A^*$  e per ogni  $a \in A$  denotiamo con  $|w|_a$  il numero di caratteri  $a$  che si incontrano scorrendo  $w$ . Evidentemente  $|w| = \sum_{a \in A} |w|_a$ .

Più in generale per ogni  $B \subseteq A$  definiamo  $|w|_B := \sum_{a \in B} |w|_a$ .

Definiamo **alfabeto minimo della parola**  $w$ , e denotiamo con  $w^{\text{minalf}}$ , il più ridotto sottoinsieme di  $A$  tale che  $w \in (w^{\text{minalf}})^*$ .

Evidentemente  $a \in w^{\text{minalf}} \iff |w|_a > 0$  e  $\mu^{\text{minalf}} = \emptyset$ .

**C10a.10** Si dice **linguaggio [formale] sull'alfabeto**  $A$  ogni sottoinsieme di  $A^*$ .

Denotiamo con  $\text{Lng}_A$  la loro collezione, cioè definiamo

$$\text{Lng}_A := \mathfrak{P}(A^*) .$$

Nello studio dei linguaggi i casi relativi ad alfabeti di una sola lettera sono molto particolari. In effetti in questo caso si hanno monoidi liberi come

$$\{ \mid \}^* = \{ \mu, \mid, \mid\mid, \mid\mid\mid, \mid\mid\mid\mid, \mid\mid\mid\mid\mid, \mid\mid\mid\mid\mid\mid, \mid\mid\mid\mid\mid\mid\mid, \dots \} .$$

Un tale insieme è evidentemente isomorfo ad  $\mathbb{N}$  e fornisce una cosiddetta **rappresentazione unadica degli interi naturali** si tratta della rappresentazione che, sostanzialmente, è stata adottata da tutti i popoli primitivi.

Va tuttavia segnalata la eccezione dello zero, entità non associabile ad alcun oggetto tangibile e quindi largamente adottato dopo un rilevante travaglio.

Secondo questa rappresentazione un intero  $n > 0$  è rappresentato da  $n$  repliche affiancate di un determinato segno (che possono rappresentare  $n$  tacche,  $n$  sassolini (calcoli),  $n$  conchiglie,  $n$  nodi su una corda, ...).

La somma di due interi naturali è ottenuta, ben comprensibilmente, accostando le sequenze di segni esprimenti i due addendi e corrisponde nei modelli concreti all'accostamento di due sequenze di tacche,

all'aggregazione di due raggruppamenti di sassolini o di conchiglie, al congiungimento di due corde che portano nodi, ... .

Per alfabeti  $A$  con due o più elementi,  $A^*$  ha una struttura molto più complessa e ricca di applicazioni di  $\mathbb{N}$ .

È quindi ragionevole pensare che i numeri naturali servano per esprimere efficacemente caratteristiche essenziali dei raggruppamenti di oggetti più o meno tangibili, ma non siano in grado di rappresentare le caratteristiche dei raggruppamenti di oggetti tra i quali sussistono rilevanti correlazioni, cioè raggruppamenti dotati di qualche complessità.

**C10a.11** Stringhe e linguaggi, nonostante la semplicità della loro definizione, costituiscono strumenti utili in moltissimi settori e vengono studiati da svariati punti di vista.

In particolare i linguaggi formali sono in grado di fornire tutte le espressioni per le entità che possono essere oggetto di studi razionali e attendibili.

Alcune proprietà dei linguaggi si ricavano con metodi algebrici che, come vedremo tra breve, fanno riferimento a semigruppamenti [B41a], semianelli [B41c] e algebre di Kleene [C32].

Questi metodi però hanno limiti piuttosto marcati: in effetti le strutture algebriche citate si basano su assiomi piuttosto deboli che portano ad algoritmi raramente stringenti quanto quelli sui numeri interi e sulle funzioni numeriche.

In effetti per disporre di linguaggi efficacemente utilizzabili si rendono necessari altri approcci di tipo costruttivo: nei capitoli C12 e C14 vedremo due classi di strumenti che consentono di individuare e manipolare linguaggi relativamente semplici, i riconoscitori e le grammatiche.

Più avanti vedremo alcuni strumenti più prestanti come le macchine di Turing, e i sistemi di riscrittura; di questi strumenti sono stati proposti numerosi arricchimenti in grado di affrontare classi mediamente estese di problemi.

Molte proprietà dei linguaggi richiedono studi di proprietà combinatorie nei quali si devono analizzare congiuntamente biezioni con configurazioni discrete di vario genere (in particolare geometrico), meccanismi enumerativi, proprietà algebriche e costruzioni algoritmiche.

Un genere di strumenti che fornisce metodi costruttivi che possono essere computazionalmente stringenti è costituito dalle serie di potenze di variabili noncommutative.

**C10a.12** Per quanto riguarda le applicazioni, occorre innanzitutto segnalare che mediante stringhe si possono esprimere, ovvero **codificare**, tutti gli oggetti studiati dalla matematica e dalle discipline che di essa si avvalgono e tutte le informazioni che si intendono sottoporre a elaborazioni automatiche.

Questo spinge a fare sistematico riferimento ai procedimenti che operano su stringhe, a cominciare dai procedimenti che servono allo studio dei fondamenti della matematica e dalla logica, fino ai procedimenti computazionali, alle tecniche per la manipolazione simbolica e all'analisi dei testi.

Molti oggetti della matematica discreta si studiano e si elaborano attraverso loro codifiche e queste rivestono un ruolo importante in vari punti della molto articolata rete di collegamenti che sussiste tra le numerose configurazioni discrete che vengono studiate anche per applicazioni di notevole rilevanza pratica.

I linguaggi formali trovano applicazioni importanti e dirette con i linguaggi di programmazione [Aho, Sethi, Ullman 1988], con i linguaggi per la comunicazione [Goldfarb 1990] e con la crittografia [Berstel, Perrin 1985].

Con essi inoltre si possono tenere sotto controllo i linguaggi naturali e vari fenomeni e processi della organizzazione delle conoscenze e della gestione della comunicazione.

Segnaliamo in particolare l'attuale importanza pratica degli algoritmi per l'analisi dei testi [Crochemore, Rytter 1994] e per la compressione dei dati [Bell, Cleary, Witten 1990].

**C10a.13** I linguaggi formali forniscono vari spunti per lo sviluppo di strumenti enumerativi.

Il più semplice proviene dalla cosiddetta **successione di enumerazione-L di un linguaggio L**, successione dei numeri delle sue stringhe raggruppate secondo le successive lunghezze:

$$\langle r \in \mathbb{N} : |L \cap A^r| \rangle.$$

Altri temi enumerativi riguardano i diversi modi di individuare le stringhe di un linguaggio quando si opera con strumenti come le grammatiche; in particolare questo conduce alla precisazione di alcuni generi di espressioni con significato ambiguo [C14].

## C10 b. operazioni su stringhe e linguaggi

**C10b.01** Nel seguito del capitolo useremo i caratteri  $A, B, \dots$  per denotare alfabeti,  $L, M, N, X, Y, \dots$  per denotare linguaggi,  $a, b, c, \dots$  per i caratteri,  $u, v, w, x, y, z$  per denotare stringhe; le lettere  $m, n, r, s$  individueranno invece interi naturali.

Si dice **giustapposizione dei linguaggi**  $L$  ed  $M$  l'insieme delle stringhe ottenibili giustapponendo una stringa di  $L$  con una di  $M$ :

$$L, M := \{w \in L, x \in M : |w, x)\}.$$

Dunque la giustapposizione di linguaggi è l'estensione booleana della giustapposizione di stringhe.

Di un linguaggio si possono quindi considerare tutte le potenze positive:

$$L^1 := L, \quad L^2 := L, L, \quad \dots L^r := L^{r-1}, L.$$

Inoltre risulta utile definire la sua potenza nulla

$$L^0 := \{\mu\}.$$

**C10b.02** Si dice **chiusura per giustapposizione del linguaggio** o **chiusura-cross del linguaggio**  $L$  l'unione delle sue potenze positive; la denoteremo con

$$L^+ := \bigcup_{r=1}^{\infty} L^r.$$

Chiaramente  $L^+ = L[ , ]$ .

Si dice invece **chiusura-star del linguaggio**  $L$  l'unione delle sue potenze nonnegative; la denoteremo con

$$L^* := \bigcup_{r=0}^{\infty} L^r.$$

Si provano facilmente relazioni come le seguenti:

$$\forall r, s \in \mathbb{N} : L^r, L^s = L^s, L^r = L^{r+s}, \quad L^* = \{\mu\} \dot{\cup} L^+, \\ L^+ = L, L^* = L^*, L; \quad L^* = L^+ \iff \mu \in L^+ \iff \mu \in L.$$

**C10b.03** In molti contesti l'operazione binaria di unione di linguaggi viene denotata con il segno “+” e conseguentemente l'unione in generale viene espressa con la notazione della sommatoria: questo atteggiamento non del tutto chiaro rende più familiari alcune proprietà algebriche dei monoidi liberi di stringhe, in quanto le avvicina a proprietà degli usuali campi numerici. Questa convenzione si accorda con quella di esprimere la giustapposizione con il segno “.” e di chiamarla “prodotto”.

Si ottengono quindi le seguenti uguaglianze:

$$L + (M + N) = (L + M) + N, \quad L \cdot (M \cdot N) = (L \cdot M) \cdot N \text{ (associatività)}, \\ (L + M) \cdot N = L \cdot N + M \cdot N, \quad L \cdot (M + N) = L \cdot M + L \cdot N \text{ (distributività)}.$$

Conviene segnalare esplicitamente che, contrariamente a quanto accade per i campi, si ha l'idempotenza della somma:

$$L + L = L;$$

quindi non si può disporre di una operazione inversa della somma.

Va anche segnalato che, in accordo con le suddette convenzioni “algebrizzanti”, talora viene usato il segno “−” per denotare la differenza insiemistica di linguaggi, cioè che si considera  $L - M$  equivalente a  $L \setminus M$  e che talora invece di  $\mu$  si usa  $\mathbf{1}$ .

Mentre ovviamente sopra un alfabeto di un carattere la giustapposizione di linguaggi è commutativa, sopra un alfabeto di due o più caratteri per gran parte delle coppie di linguaggi  $\langle L, M \rangle$  si riscontra la noncommutatività, cioè accade che  $L \cdot M \neq M \cdot L$ .

Per i linguaggi  $\emptyset$  e  $\{\mu\}$  si trovano le seguenti relazioni:

$$\begin{aligned} L + \emptyset = \emptyset + L = L \quad , \quad L \cdot \emptyset = \emptyset \cdot L = \emptyset \quad , \quad L \cdot \mu = \mu \cdot L = L \quad , \\ \emptyset^+ = \emptyset \quad , \quad \emptyset^* = \mu \quad , \quad \mu^+ = \mu^* = \mu \quad . \end{aligned}$$

L'insieme dei linguaggi sopra un alfabeto  $A$  costituisce quindi un semianello avente  $\emptyset$  come zero e  $\{\mu\}$  come unità; questo semianello è noncommutativo sse  $A$  ha più di un carattere [T15h02].

Nell'ultima delle precedenti relazioni abbiamo adottato un'altra semplificazione piuttosto usuale evitando le parentesi graffe intorno a  $\mu$ , cioè identificando una stringa con il linguaggio costituito da questo solo elemento.

**C10b.04** Utilizzando i segni “+” e “.” le espressioni per i linguaggi finiti assumono aspetto di espressioni polinomiali: ad esempio invece di  $\{a^3, b^3\}$  si scrive  $a^3 + b^3$ , cioè una espressione utilizzata anche per oggetti molto familiari come i polinomi.

Questo agevola la scrittura di formule come la seguente:

$$(a + b^2 + b^2 a) \cdot (\mu + a) = a + b^2 + b^2 a + a^2 + b^2 a^2 .$$

In queste espressioni l'idempotenza della somma non consente coefficienti superiori a 1 e ovviamente non hanno senso i coefficienti negativi; i soli coefficienti delle espressioni monomiali possono essere lo zero e l'unità, ovvero  $\emptyset$  e  $\mu$ ; questi due elementi costituiscono il terreno del semianello isomorfo a quello dei bits  $\mathbb{B}$ .

In queste espressioni i caratteri svolgono un ruolo che li fa considerare delle ‘variabili noncommutative’ sopra un semianello.

Dunque le stringhe e linguaggi sono retti da formule contenenti espressioni polinomiali in variabili noncommutative nelle quali non compare alcun coefficiente.

Coerentemente un linguaggio infinito  $L \subseteq A^*$  si può considerare come serie nelle variabili formali noncommutative costituenti  $A$  sul semianello avente come terreno  $\{\emptyset, \mu\}$ .

Questo conduce a scritture come la seguente:

$$L = \sum_{w \in L} w = \sum_{\mathbf{a}_{j_1} \dots \mathbf{a}_{j_r} \in A^*} \mathbf{a}_{j_1} \dots \mathbf{a}_{j_r} \langle \mathbf{a}_{j_1} \dots \mathbf{a}_{j_r} | L \rangle ,$$

nella quale si sono introdotti i coefficienti (alla Dirac) con valori binari

$$\langle \mathbf{a}_{j_1} \dots \mathbf{a}_{j_r} | L \rangle := \begin{cases} \mu & \text{sse } \mathbf{a}_{j_1} \dots \mathbf{a}_{j_r} \in L, \\ \emptyset & \text{altrimenti.} \end{cases}$$

In particolare si hanno la serie delle potenze della stringa  $w$  e la serie delle potenze del linguaggio  $L$ :

$$w^* = \mu + w + w^2 + \dots = \sum_{r=0}^{\infty} w^r \quad , \quad L^* = \mu + L + L^2 + \dots = \sum_{r=0}^{\infty} L^r .$$

Quindi per affermare che una stringa  $z$  è una qualche potenza della  $w$  si scrive  $z \in w^*$ , mentre per enunciare che tutte le stringhe del linguaggio  $M$  si possono esprimere come prodotti di stringhe appartenenti ad  $L$  si scrive  $M \subseteq L^*$ .

**C10b.05** Per le operazioni definite sui linguaggi si possono dimostrare senza difficoltà varie altre relazioni.

**Prop.**  $\forall L, M \subseteq A^*$  :

$$\begin{aligned} (L + M)^* &= (L^* \cdot M)^* \cdot L^* = (M^* \cdot L)^* \cdot M^* \quad , \quad (L \cdot M)^* = \mu + L \cdot (M \cdot L)^* \cdot M \quad , \\ L^* \cdot L^* &= L^* \quad , \quad L^+ \cdot L^+ = L \cdot L^+ = L^2 \cdot L^* \quad , \quad L^* \cdot L^+ = L^+ \cdot L^* = L^+ \quad , \\ (L^*)^* &= L^* \quad , \quad (L^+)^+ = L^+ \quad , \quad (L^+)^* = (L^*)^+ = L^* \quad ; \quad (\mu + L)^+ = (\mu + L)^* = L^* \quad , \\ \forall r \in \mathbb{N} : (L^+)^r &= L^r \cdot L^* \quad , \quad \forall r \in \mathbb{N} : L^* = (L^r)^* \cdot L^{[r]} \quad \blacksquare \end{aligned}$$

**C10b.06** Le uguaglianze trovate che coinvolgono gli operatori unari  $^+$  e  $^*$  inducono a riferire le costruzioni insiemistiche e basate sulla giustapposizione dei linguaggi a una struttura algebrica più ricca della struttura di semianello.

Seguendo *John Horton Conway (1971)* chiamiamo **algebra di Kleene classica** una struttura della forma  $\langle \mathbf{K}, \oplus, \odot, *, \mathbf{0}, \mathbf{1} \rangle$  nella quale

- $\langle \mathbf{K}, \oplus, \odot, \mathbf{0} \rangle$  è un semianello ;
- l'operazione binaria “ $\oplus$ ” è idempotente:  $(\forall L \in \mathbf{K} : L \oplus L = L)$  ;
- “ $\mathbf{1}$ ” è l'unità per l'operazione “ $\odot$ ” ,
- “ $*$ ” è un'operazione unaria,

per la quale valgono le seguenti proprietà per arbitrari  $L, M \in \mathbf{K}$ :

$$\begin{aligned} (L \oplus M)^* &= (L^* \odot M)^* \odot L^* \quad , \quad (L \odot M)^* = \mathbf{1} \oplus L \odot (M \odot L)^* \odot M \quad , \\ L^* \odot L^* &= L^* \end{aligned}$$

e dove, posto  $L^\oplus := L \odot L^*$  , vale la successione di uguaglianze

$$\forall r \in \mathbb{N} : L^* = (L^r)^* \odot (\mathbf{1} \oplus L \oplus \dots \oplus L^{r-1}) \quad .$$

Oltre la ridottissima struttura con il terreno ridotto a un solo elemento che fa da zero e da unità, la più semplice algebra di Kleene classica che chiamiamo **algebra di Kleene booleana** è

$$\langle \{\emptyset, \mu\}, +, \cdot, *, \emptyset, \mu \rangle$$

Molto più articolata è l'**algebra di Kleene dei linguaggi** sull'alfabeto  $A$ :

$$\langle \mathbf{Lng}_A, +, \cdot, *, \emptyset, \mu \rangle \quad .$$

La struttura algebrica precedente e altre strutture a lei vicine verranno trattate con una certa ampiezza in C32.

**C10b.07** Una operazione unaria su stringhe e linguaggi che conduce a una utile simmetria negli insiemi di stringhe e linguaggi è la **riflessione**.

Si dice **stringa riflessa** della stringa  $w =_A a_{j_1} a_{j_2} \dots a_{j_r}$  e si denota con  $w^\leftarrow$  o  $w^\top$  la stringa ottenuta “leggendo i caratteri della  $w$  da destra a sinistra”, cioè la stringa  $w^\leftarrow := a_{j_r} \dots a_{j_2} a_{j_1}$  .

Ad esempio:  $(ROMA)^\leftarrow = AMOR$  .

Per la giustapposizione di due stringhe si constata che  $(w \cdot x)^\leftarrow = x^\leftarrow \cdot w^\leftarrow$  ; questo si esprime in termini algebrici dicendo che l'operazione unaria riflessione è un **antiisomorfismo dei monoidi liberi**.

Si dice **riflesso del linguaggio**  $L$ , e si scrive  $L^\leftarrow$ , l'insieme delle stringhe riflesse delle stringhe di  $L$ :

$$\forall L \in \mathbf{Lng} : L^\leftarrow := \{w \in L : | w^\leftarrow\} \quad .$$

Evidentemente la riflessione dei linguaggi è l'estensione booleana della riflessione tra stringhe; inoltre essa è un'involuzione per la collezione dei linguaggi sopra un qualsiasi alfabeto:

$$\forall L \in \mathbf{Lng} : (L^{\leftarrow})^{\leftarrow} = L .$$

Tra le stringhe si distinguono quelle invarianti per riflessione, stringhe dette **stringhe palindrome** : una palindroma di lunghezza pari è ANNA ed una di lunghezza dispari è ANILINA. un'altra palindroma famosa è SATOR AREPO TENET OPERA ROTAS .

Convieni considerare che anche la stringa muta sia una palindroma.

Un linguaggio P si dice **linguaggio palindromo** sse  $w \in P \implies w^{\leftarrow} \in P$ .

Denotiamo con  $\mathbf{LngPIndrm}_A$  la collezione dei linguaggi palindromi su A e con  $\mathbf{Palnd}_A$  l'insieme delle stringhe palindrome su A.

(1) **Prop.:** La collezione dei linguaggi costituiti da sole stringhe palindrome su A è contenuta propriamente nella collezione dei linguaggi palindromi su A, ossia  $\mathfrak{P}(\mathbf{Palnd}_A) \subset \mathbf{LngPIndrm}_A$  .

**Dim.:** Chiaramente  $\forall P \subseteq \mathbf{Palnd}_A : P \in \mathbf{LngPIndrm}_A$  ; mentre per l'inclusione stretta basta osservare che  $aab + baa \in \mathbf{LngPIndrm}_{\{a,b\}}$  mentre  $(aab + baa) \cap \mathbf{Palnd}_{\{a,b\}} = \emptyset$  .

**C10b.08 (1) Eserc.** Verificare che l'insieme delle palindrome su A si può esprimere come:

$$\mathbf{Palnd}_A := \sum_{w \in A^*} w w^{\leftarrow} + \sum_{w \in A^*} w A w^{\leftarrow}$$

e che, se  $n := |A|$ , la successione di enumerazione-L di tale linguaggio è:

$$\left| \begin{array}{cccccccc} 0 & 1 & 2 & 3 & 4 & 5 & 6 & \dots & r & \dots \\ \downarrow & n & n & n^2 & n^2 & n^3 & n^3 & \dots & n^{\lceil r/2 \rceil} & \dots \end{array} \right| .$$

(2) **Eserc.** Costatare che  $\forall y \in \mathbf{Palnd}_A : y^2, y^3, \dots, y^+, y^* \in \mathbf{LngPIndrm}_A$  .

(3) **Eserc.** Costatare che  $\mathbf{LngPIndrm}_A \in [ +, \leftarrow ]$

(4) **Eserc.** Mostrare controesempi che portano alla relazione  $\mathbf{LngPIndrm}_A \notin [ \cdot ]$  e constatare che si può solo affermare  $\lceil Y, Z \in \mathbf{LngPIndrm}_A \implies Y \cdot Z + Z \cdot Y \in \mathbf{LngPIndrm}_A \rceil$  .

(5) **Eserc.** Costatare che  $Y \subseteq \mathbf{Palnd}_A \not\Rightarrow Y^+, Y^* \in \mathbf{LngPIndrm}_A$  e quindi  $Z \in \mathbf{LngPIndrm}_A \not\Rightarrow Z^+, Z^* \in \mathbf{LngPIndrm}_A$  .

(6) **Eserc.** Costatare che  $\forall L \in \mathbf{Lng}_A : L + L^{\leftarrow} \in \mathbf{LngPIndrm}_A \wedge (L + L^{\leftarrow})^+ \in \mathbf{LngPIndrm}_A$

**C10b.09** La riflessione e la giustapposizione sono state definite sui linguaggi attraverso la semplice estensione booleana delle corrispondenti operazioni definite sulle stringhe. Analoghe **estensioni booleane** si possono effettuare per tutte le operazioni unarie, binarie e infinitarie definite sugli elementi dei linguaggi e di una qualsiasi struttura algebrica per renderle applicabili ai sottoinsiemi (del terreno) della struttura stessa.

Se si adotta la stessa notazione per una operazione e per la sua estensione booleana in genere non si generano confusioni difficili da chiarire osservando il contesto.

Questo si riscontra in particolare se si semplifica il simbolo  $\odot^{be}$  con  $\odot$  e il simbolo  $+^{be}$  con  $+$ .

Esaminando le estensioni booleane si giustificano facilmente le uguaglianze che seguono, uguaglianze che adottano la semplificazione suddetta e che aiutano a chiarire le relazioni tra la riflessione e le altre operazioni sui linguaggi:

$$\begin{aligned} (L + M)^{\leftarrow} &= L^{\leftarrow} + M^{\leftarrow} & , & & (L \cdot M)^{\leftarrow} &= M^{\leftarrow} \cdot L^{\leftarrow} & , \\ (L^+)^{\leftarrow} &= (L^{\leftarrow})^+ & , & & (L^*)^{\leftarrow} &= (L^{\leftarrow})^* & . \end{aligned}$$

Ricordiamo inoltre che la riflessione tra stringhe e linguaggi è alla base della dualità-LR.

**C10b.10** Per talune attività, dato un linguaggio  $L$  serve individuare il suo cosiddetto sottolinguaggio palindromo

$$(1) \quad L \cap \mathit{Palnd}.$$

In altre circostanze si chiede il suo sottolinguaggio antipalindromo  $M \subseteq L$  tale che  $M \cap M^{\leftarrow} = \emptyset$ , cioè l'insieme ottenuto eliminando da  $L$  le sue stringhe palindrome. Questo linguaggio è fornito dalle espressioni equivalenti

$$(2) \quad L \setminus \mathit{Palnd} = L \cap \mathit{Palnd} = L \setminus (L \cap L^{\leftarrow}).$$

Occorre osservare che le espressioni proposte (1) e (2) quando si devono trovare caratteristiche utilizzabili dei due sottolinguaggi, servono solo ad impostare i due problemi effettivi, ma la loro soluzione deve essere trovata tenendo conto delle specifiche caratteristiche del linguaggio da ridurre e delle esigenze che si vogliono soddisfare con il sottolinguaggio risultato.

**C10b.11** Nella lingua italiana, ma non solo, si trovano molte parole ottenute dalla giustapposizione di due sillabe con il frequente raddoppio della iniziale della seconda:

babà, bebè, cacca, cocco, lalla, mamma, nanna, nonno, papa, pappà, popo, pupu, tata, tette, ...

Molte di queste parole sono tra le prime che i piccoli sono in grado di controllare, grazie a una loro sostanziale semplicità costruttiva; la presenza di una doppia ha l'effetto di una pausa che facilita ulteriormente in controllo della parola.

Un altro schema che raccoglie parole di facile controllo consiste in un tetragramma palindromo preceduto da una consonante: cassa, nassa, sonno, zappa, ... .

Segnaliamo anche le giustapposizioni di due palindrome: esse godono di due proprietà di simmetria abbastanza semplici, anche se meno nette della proprietà che caratterizza le palindrome.

Se  $u$  e  $v$  sono due palindrome, allora  $(u_1 v)^{\leftarrow} = v_1 u$ .

**C10b.12** Buona parte dei giochi enigmistici consistono nel presentare un testo, tipicamente un testo poetico, nel quale sono individuabili due o più parole incognite che si chiede di individuare. Queste parole vengono rappresentate da schemi convenzionali (xxxxxx, xxxyyyxxx, xxyxxxzxx, ...), oppure che sono alluse da immagini espresse in versi collegati da rime e somiglianze ritmiche, spesso facendo intervenire termini ambigui per omonimia, sinonimia e polisemia.

Le parole incognite risultano raggiungibili attraverso collegamenti esprimibili in termini di operazioni su stringhe come giustapposizione, riflessione o permutazione.

Il solutore deve quindi ricostruire percorsi basati su collegamenti criptici che possono essere sia formali che dovuti a significati non immediatamente riconoscibili.

Si parla di “anagramma” se si devono indovinare una parola o un frase e una sua permutata (*trilussa* e *salustri*), (*arrosti* con *artrosi*).

Più in particolare si parla di “scambio” quando le due parole differiscono solo per il cambiamento di posizione di due occorrenze di lettere (*creole* con *ercole*).

Si chiama “scambio sillabico” uno scambiano di due sillabe.

Si chiama “bifronte” il gioco nel quale si devono indovinare una parola e la sua riflessa (*acetone* insieme a *enoteca*).

Si parla di “cambio” o di “metatesi” quando si modifica una sola lettera (*malta*, *malva*).

Si ha una “zeppa” se da una prima parola si ottiene una seconda mediante inserimento di una lettera (da *casa* a *casta*) o di più lettere (da *lino* a *liocorno*).

Viceversa viene detto “scarto” il gioco che vede una prima parola e una seconda ottenuta eliminando dalla prima una lettera (da *geloso* a *gelso*).

Loro prevedibili varianti sono la “zeppa sillabica” (*ditale* con *digitale*, *scorno*, *scorfano* e *matto* con *mattone*, che più in particolare è un esempio di falso accrescitivo) e lo “scarto sillabico” (*anonimo*, *animo*).

Una “sciarada” richiede di trovare due parole che giustapposte forniscono una terza parola o frase incognita (*ciocco* e *latino* danno *cioccolatino*)

## C10 c. morfismi tra linguaggi

**C10c.01** Ricordiamo che in algebra un sottoinsieme  $T$  di un generico semigruppò  $S$  si dice: **ideale da sinistra** sse  $ST = T$ , cioè sse è un sottoinsieme invariante per traslazioni da sinistra; **ideale da destra** sse  $TS = T$ , cioè sse è un sottoinsieme invariante per traslazioni da destra; **ideale bilatero** sse  $STS = T$ , cioè sse è un sottoinsieme invariante per traslazioni da destra e da sinistra; **sottosemigruppò** sse  $TT \subseteq T$ .

Evidentemente gli ideali bilateri sono anche ideali da sinistra e ideali da destra.

Inoltre i sottosemigruppò sono anche ideali da sinistra e ideali da destra (nonché ideali bilateri).

Un sottoinsieme  $N$  di un monoide  $M$  avente  $\mathbf{1}$  come unità si dice **sottomonoido** di  $M$  sse  $\mathbf{1} \in N$  e  $NN \subseteq N$ , ovvero sse  $NN = N$ .

Esaminiamo questi tipi di sottostrutture nel caso del monoide libero  $A^*$ .

**C10c.02 Prop.** Consideriamo i linguaggi sopra un alfabeto  $A$ :

I linguaggi ideali da sinistra sono esattamente i linguaggi della forma  $A^*M$  per qualche  $M \in \mathbf{Lng}_A$ ; in altre parole: la collezione dei linguaggi ideali a sinistra è  $\{M \in \mathbf{Lng}_A : A^*M\}$ .

In accordo con la dualità-LR, l'insieme dei linguaggi ideali da destra è  $\{M \in \mathbf{Lng}_A : MA^*\}$ .

L'insieme dei linguaggi ideali bilateri è  $\{M \in \mathbf{Lng}_A : A^*MA^*\}$ .

Si noti che  $A^*M \cap MA^* = A^*MA^*$ .

I linguaggi sottosemigruppò hanno forma  $T = M^+$  per qualsiasi  $M \in \mathbf{Lng}_A$ , cioè sono i linguaggi tali che  $T^+ = T$ .

I linguaggi sottomonoidi hanno forma  $T = M^*$ , cioè sono i linguaggi tali che  $T^* = T$ ; equivalentemente possiamo affermare che l'insieme dei sottomonoidi di  $A^*$  è  $\{L \subseteq A^* : L^*\}$ .

Relativamente all'alfabeto  $A$ , si dicono **linguaggi cometa** i linguaggi della forma  $A^*N$ ; sono detti invece **linguaggi anticometa** quelli della forma  $MA^*$ ; si chiamano invece **linguaggi bicometa** i linguaggi della forma  $A^*NA^* = A^*N \cap NA^*$ .

Si constata che i riflessi delle comete sono le anticomete, mentre le collezioni delle bicomete, dei sottosemigruppò e dei sottomonoidi sono invarianti per riflessione.

Inoltre possiamo affermare che un linguaggio bicometa della forma  $A^*NA^*$  è un linguaggio palindromo sse  $N$  è palindromo.

Possiamo inoltre affermare che rispetto alla riflessione le collezioni degli ideali da sinistra e degli ideali da destra sono mutuamente duali-LR, mentre le collezioni dei linguaggi ideali bilateri, dei linguaggi sottosemigruppò e dei linguaggi sottomonoidi sono autoduali-LR.

**(1) Eserc.** Costatare che le collezioni dei linguaggi cometa, quella dei linguaggi anticometa e quella dei linguaggi bicometa sono collezioni (numerabili) invarianti per somma e giustapposizione e di conseguenza per chiusura-cross e per sommatoria numerabile, ovvero che appartengono alla collezione  $[+, \cdot, \oplus, \sum^{+\infty}]$ .

**C10c.03** Abbiamo visto che  $\langle L^+, \cdot \rangle \in \mathbf{Sgrp}$  e che  $\langle L^*, \cdot, \mu \rangle \in \mathbf{Mnd}$ ; possiamo quindi affermare che  $L^+$  è il sottosemigruppò di  $A^+$  generato da  $L$ , mentre  $L^*$  è il sottomonoido di  $A^+$  generato da  $L$ .

In accordo con questa terminologia algebrica, possiamo affermare che semigruppò liberi e monoidi liberi sono, risp., i semigruppò e i monoidi generati dagli alfabeti.

Dato un sottomonoido  $M$  di  $A^*$ , è spesso utile individuare un suo insieme di generatori, cioè un linguaggio contenuto in  $A^*$  che generi  $M$ .

Evidentemente sono più interessanti i sistemi di generatori più ridotti, cioè minimali per l'inclusione.

**(1) Prop.:** Ogni sottomonoido  $M$  di  $A^*$  ha un unico insieme di generatori minimale per inclusione individuabile con l'espressione di impostazione:

$$(M \setminus \mu) \setminus (M \setminus \mu)^2.$$

**Dim.:** Ogni insieme di generatori di  $M$  deve essere contenuto in  $M \setminus \mu$  e per essere minimale non può contenere alcuna stringa in  $(M \setminus \mu)^2$ , cioè è contenuto nel linguaggio definito dall'espressione.

Viceversa ogni stringa di  $M$  è ottenibile come prodotto di elementi del linguaggio sopra espresso ■

Convieni ribadire che l'espressione precedente (che spesso si trova scritta  $(M - \mu) - (M - \mu)^2$ ) fornisce solo una indicazione genericamente formale per la ricerca dell'insieme di generatori minimale di  $M$ .

Se si rende necessario controllare efficacemente l'insieme generatore minimale di uno specifico  $M \leq_{Mnd} A^*$  (in particolare se occorre decidere se tale insieme è finito o infinito) occorre individuare un procedimento effettivo che si basa sulle proprietà specifiche che determinano  $M$ .

Le stesse considerazioni valgono per le espressioni **b10(1)** e **(2)**.

**C10c.04** In generale si pone il problema della utilizzabilità delle espressioni matematiche, ossia delle formule che esprimono risultati matematici.

Alcune delle espressioni come la **(1)** e le due sopraccitate servono solo ad inquadrare un problema che dovrà essere approfondito tenendo conto di proprietà specifiche.

All'opposto talune espressioni, in particolare le uguaglianze che si incontrano nell'algebra elementare e nell'analisi numerica, contengono tutte le indicazioni che consentono calcoli effettivi e lasciano aperti solo questioni riguardanti le tolleranze consentite alle approssimazioni.

Vi sono poi situazioni intermedie rappresentate da espressioni contenenti operazioni ben definite, cioè calcolabili con algoritmi consolidati, e operazioni definite in modo generico, vaghe, che richiedono ulteriori approfondimenti che coinvolgono operandi che non si sanno completamente dominare.

Per gli approfondimenti che molte formule lasciano in sospeso riveste notevole importanza il livello di generalità consentito per gli argomenti che occorrono nelle espressioni. Tipicamente contano gli operandi e i domini delle funzioni e delle relazioni.

La portata di una formula cresce al crescere dei domini di variabilità dei suoi argomenti. Questa tendenza si può trovare in conflitto con la utilizzabilità effettiva dell'espressione.

Una espressione riguardante domini più estesi tende ad essere meno utilizzabile di una versione della stessa formula riguardante domini meno estesi.

I compromessi tra queste esigenze conflittuali sono oggetto di studi di notevole importanza. Queste tematiche sono definite in modo particolarmente chiaro quando espresse in termini di digrafi e di linguaggi formali.

**C10c.05** Ricordiamo che si dice **morfismo** dal monoide  $\langle M, \cdot, 1 \rangle$  nel monoide  $\langle N, \odot, 1 \rangle$  ogni applicazione  $\psi \in [M \mapsto N]$  tale che  $\forall u, v \in M : \psi(u \cdot v) = \psi(u) \odot \psi(v)$  e  $\psi(1) = 1$ .

**(1) Prop.:** Consideriamo un insieme  $G$ , un monoide  $M$ , una funzione  $f \in [G \mapsto M]$  e l'iniezione naturale  $i := \text{natinj}(G, G^*) = [g \in G \mapsto G^*]$ .

Esiste un unico morfismo di monoidi  $\phi \in [G^* \mapsto M]$  tale che  $f = \text{natinj}(G, G^*) \circ_{lr} \phi$ , cioè tale da rendere commutativo il diagramma funzionale

//input pC10c05 ■

Questa proprietà di  $G^*$ , chiamata **proprietà universale** equivale ad affermare che  $G^*$  sia un **monoide libero**.

La situazione della proposizione (1) nel caso  $G$  sia un alfabeto finito  $A$  conduce al morfismo  $\psi$  dal monoide libero  $A^*$  in un monoide  $\langle M, \odot, \mathbf{1} \rangle$  che risulta definito quando si forniscono i trasformati  $\psi(\mathbf{a}_i)$  per tutti i caratteri  $\mathbf{a}_i \in A$ , cioè per tutti gli elementi generatori del dominio: infatti per ogni  $\mathbf{a}_{j_1}, \mathbf{a}_{j_2}, \dots, \mathbf{a}_{j_r} \in A^*$  si ottiene

$$\psi(\mathbf{a}_{j_1} \mathbf{a}_{j_2} \dots \mathbf{a}_{j_r}) = \psi(\mathbf{a}_{j_1}) \odot \mathbf{a}_{j_2} \odot \dots \odot \psi(\mathbf{a}_{j_r}) .$$

**C10c.06** I morfismi tra due monoidi liberi  $A^*$  e  $B^*$  sono anche detti, motivatamente, **sostituzioni di caratteri con stringhe**.

Ciascuno di tali morfismi  $\psi$  è definito dalla funzione  $\gamma = \{ \mathbf{a}_i \in A \mapsto \psi(\mathbf{a}_i) \}$  e gli si può assegnare la forma  $\psi = \gamma^\otimes = \{ \mu \mapsto \mu \} \cup \gamma \cup \gamma \circ \gamma \cup \gamma^{\circ 3} \cup \dots$ .

Sono interessanti alcuni sottoinsiemi specifici di un insieme di morfismi  $\mathbf{M} := \{ A^* \mapsto B^* \}$ .

Un morfismo  $\psi \in \mathbf{M}$  si dice **cancellazione di caratteri** sse  $\psi(\mathbf{a}_i) = \mu$  per qualche  $\mathbf{a}_i \in A$ .

SLa funzione  $\psi \in \mathbf{M}$  si dice **funzione che rispetta la lunghezza** sse per ogni  $\mathbf{a}_i \in A$  si ha  $\psi(\mathbf{a}_i) \in B$ .

Le cancellazioni di caratteri non rispettano la lunghezza ■

Un morfismo  $\gamma^\otimes =: \psi$  va chiamato **isomorfismo** sse  $\gamma \in \{ A \leftrightarrow B \}$ , ovvero sse

$$\forall \mathbf{a}_i, \mathbf{a}_j \in A \quad \psi(\mathbf{a}_i) = \psi(\mathbf{a}_j) \implies \mathbf{a}_i = \mathbf{a}_j .$$

Tale trasformazione rispetta la lunghezza.

Si chiamano **morfismi invertibili** i morfismi del genere  $A^* \leftrightarrow B^*$ ; si tratta di trasformazioni invertibili e quindi tali da consentire, per ogni  $w \in A^*$ , la ricostruzione di  $w$  a partire da  $\psi(w)$  come  $w = (\psi^{-1})(\psi(w))$ .

Queste trasformazioni sono l'oggetto della teoria dei codici, una disciplina, oggi assai sviluppata e sostegno di importanti applicazioni [C65, C66].

**C10c.07** Per estensione booleana le sostituzioni si applicano anche ai linguaggi definendo:

$$(1) \quad \forall L \in \mathbf{Lng}_A \quad \psi(L) := \sum_{\mathbf{a}_{j_1} \mathbf{a}_{j_2} \dots \mathbf{a}_{j_r} \in L} \psi(\mathbf{a}_{j_1}) \cdot \psi(\mathbf{a}_{j_2}) \cdot \dots \cdot \psi(\mathbf{a}_{j_r}) .$$

Questa espressione si può utilizzare per definire più in generale l'applicazione a un linguaggio su un alfabeto  $A$  di una funzione del genere  $\psi \in \{ A \mapsto \mathbf{Lng}_B \}$  la quale a ogni carattere di  $A$  fa corrispondere un linguaggio su  $B$ .

Questa funzione si dice **sostituzione di caratteri con linguaggi**; essa non costituisce un morfismo tra monoidi liberi ma un morfismo tra algebre di Kleene di linguaggi. Infatti si dimostra facilmente che:

$$(2) \quad \psi(L + M) = \psi(L) + \psi(M) \quad , \quad \psi(L \cdot M) = \psi(L) \cdot \psi(M) \quad , \quad \psi(L^*) = (\psi(L))^* .$$

Ad esempio., se  $\psi(\mathbf{a}) := c d^+ c$  e  $\psi(\mathbf{b}) := c d c$  si ha:

$$\psi(\mu + \mathbf{a}^2 + \mathbf{b} \mathbf{a}^2 \mathbf{b}) = \mu + c d^+ c^2 d^+ c + c d c^2 d^+ c^2 d^+ c^2 d c .$$

**C10c.08** Gli isomorfismi del genere  $\{ A^* \leftrightarrow \{0, 1\}^* \}$  sono chiamati anche **codifiche binarie**. Essi consentono di rappresentare mediante sequenze di bits i più svariati tipi di stringhe e quindi i più diversi tipi di informazioni mantenendo la possibilità di risalire alla loro forma originaria.

Le rappresentazioni binarie delle informazioni consentono di registrarle, analizzarle e modificarle servendosi degli odierni strumenti dell'informatica e delle telecomunicazioni e quindi con grandi vantaggi in termini di velocità, di versatilità e di componibilità delle elaborazioni e con la conseguente possibilità di raccogliere, strutturare, riutilizzare e far circolare grandi collezioni di dati [v.a **big data** (we)].

In particolare si hanno le codifiche binarie delle decine di caratteri utilizzati correntemente con strumenti elettronici seguendo il codice **ASCII** (wi), e le codifiche binarie delle decine di migliaia di segni degli alfabeti delle lingue di qualche importanza seguendo il codice **Unicode** (wi).

I procedimenti e i metodi che consentono di codificare e decodificare le informazioni afferiscono alla sopraccitata teoria dei codici.

## C10 d. relazioni tra stringhe e derivazioni di linguaggi

**C10d.01** Consideriamo  $u, v, w \in A^*$ .

Si dice che  $u$  è **prefisso** o **fattore sinistro** di  $w$ , e si scrive  $u \preceq_p w$ , sse  $w \in uA^*$ .

Si dice che  $u$  è **suffisso**, **postfisso** o **fattore destro** di  $w$ , e si scrive  $u \preceq_s w$  sse  $w \in A^*u$ .

Si dice che  $u$  è **infisso** o **fattore** di  $w$ , e si scrive  $u \preceq_i w$ , sse  $w \in A^*uA^*$ .

Si dice che  $u$  è **sottostringa** o **sottosequenza** di  $w$ , e si scrive  $u \preceq_u w$ , sse, essendo  $u =_A a_{j_1} a_{j_2} \cdots a_{j_r}$ , si ha

$$w \in A^* a_{j_1} A^* a_{j_2} A^* \cdots A^* a_{j_r} A^* .$$

Le relazioni tra stringhe  $\preceq_p, \preceq_s, \preceq_i$  e  $\preceq_u$  sono evidentemente delle relazioni d'ordine parziale; questo ordine non è totale sse  $|A| \geq 2$ .

Accanto a queste quattro relazioni si introducono quattro funzioni del genere  $\lceil A^* \mapsto \mathbf{Lng}_A \rceil$  che ad ogni stringa  $w \in A^*$  associano:

l'insieme dei suoi prefissi:  $\mathbf{Pfx}(w) = w^{\mathbf{Pfx}}$  ;

l'insieme dei suoi suffissi:  $\mathbf{Sfx}(w) = w^{\mathbf{Sfx}}$  ;

l'insieme dei suoi infissi:  $\mathbf{Ifx}(w) = w^{\mathbf{Ifx}}$  ;

l'insieme delle sue sottosequenze (o sottotringhe):  $\mathbf{Useq}(w) = w^{\mathbf{Useq}} = \mathbf{Sbseq}(w)$ .

L'insieme dei prefissi di  $w = \mathbf{abaa}$  è  $\mathbf{Pfx}(w) = \{\mu, a, ab, aba, abaa\}$  ;

per l'insieme dei suffissi di  $w$  abbiamo  $w^{\mathbf{Sfx}} = \mathbf{Sfx}(w) = \{\mu, a, aa, baa, abaa\}$  ;

l'insieme dei suoi infissi è  $\mathbf{Ifx}(w) = \{\mu, a, b, aa, ab, ba, aba, baa, abaa\}$  ;

Per avere l'insieme delle sottotringhe di  $w$   $\mathbf{Useq}(w)$  occorre aggiungere all'insieme dei suoi infissi  $\mathbf{a}^3$ , cioè  $w^{\mathbf{Useq}} = \{\mu, a, b, aa, ab, ba, aaa, aba, baa, abaa\}$  .

**C10d.02** Sono utili anche le estensioni booleane delle 4 funzioni precedenti che noi denoteremo con gli stessi simboli; esse appartengono al genere  $\lceil \mathbf{Lng}_A \mapsto \mathbf{Lng}_A \rceil$  . Si segnala anche che queste funzioni applicate a stringhe portano a linguaggi finiti, ossia a linguaggi polinomiali.

Esse possono essere considerate endofunzioni entro  $\mathbf{Lng}_A$  e questo consente di comporre queste ed altre funzioni che trasformano linguaggi in linguaggi formali.

Ad esempio si osserva che l'insieme degli infissi di  $w$  è un linguaggio ottenibile come insieme dei prefissi dei vari suffissi di  $w$  e coincide con l'insieme dei suffissi dei vari prefissi di  $w$ , ossia

$$\forall w \in ASs^* : \mathbf{Pfx}(\mathbf{Sfx}(w)) = \mathbf{Sfx}(\mathbf{Pfx}(w)) = \mathbf{Ifx}(w) .$$

Questa doppia uguaglianza vale anche quando l'argomento è un linguaggio, ossia

$$\forall L \subseteq A^* : \mathbf{Pfx}(\mathbf{Sfx}(L)) = \mathbf{Sfx}(\mathbf{Pfx}(L)) = \mathbf{Ifx}(L) .$$

Si constata facilmente che le endofunzioni entro  $A^*$   $\mathbf{Pfx}$ ,  $\mathbf{Sfx}$ ,  $\mathbf{Ifx}$  e  $\mathbf{Useq}$  sono idempotenti.

Queste quattro endofunzioni consentono di individuare quattro collezioni di linguaggi particolari caratterizzati come loro punti fissi.

La collezione  $\lceil \mathbf{Pfx} \rceil$  dei linguaggi invarianti per  $\mathbf{Pfx}$ , linguaggi  $L$  tali che  $(\mathbf{Pfx}(L)) = L$ .

La collezione  $\lceil \mathbf{Sfx} \rceil$  dei linguaggi invarianti per  $\mathbf{Sfx}$ , linguaggi  $L$  tali che  $(\mathbf{Sfx}(L)) = L$ .

La collezione  $\lceil \mathbf{Ifx} \rceil$  dei linguaggi invarianti per  $\mathbf{Ifx}$ , linguaggi  $L$  tali che  $(\mathbf{Ifx}(L)) = L$ .

La collezione  $\lceil \mathbf{Useq} \rceil$  dei linguaggi invarianti per  $\mathbf{Useq}$ , linguaggi  $L$  tali che  $(\mathbf{Useq}(L)) = L$ .

Si osserva che i linguaggi invarianti per **Useq** sono invarianti anche per **lfx**, ma vi sono linguaggi invarianti per **lfx** che non sono invarianti per **Useq**; concisamente abbiamo

$$[ \text{Useq} ] \subset [ \text{lfx} ] .$$

Si trova inoltre che  $[ \text{lfx} ] \subset [ \text{Pfx} ]$  e dualmente-LR abbiamo  $[ \text{lfx} ] \subset [ \text{Sfx} ]$  .

**C10d.03** Munendo  $A^*$  con ciascuna delle relazioni  $\preceq_x$  per  $x \in \{p, s, i, u\}$  si ottiene un digrafo infinito graduato, i cui livelli contengono elementi delle successive potenze cartesiane  $A^r$ .

Quando  $A = \{1\}$  (ed in genere quando  $A$  è costituito da un solo carattere) si ha il semplice digrafo numerabile e funzionale

$$\mu \longrightarrow 1 \longrightarrow 11 \longrightarrow 111 \longrightarrow 1111 \longrightarrow \dots$$

isomorfo a quello di  $\mathbb{N}$  munito della relazione “immediatamente minore”  $<_1$ .

Si osserva che per ogni linguaggio  $M$  su un solo carattere, in particolare su  $1$ , abbiamo

$$\text{Pfx}(M) = \text{Sfx}(M) = \text{lfx}(M) = \text{Sbseq}(M) .$$

Inoltre se  $M$  è finito  $\text{Pfx}(M) = \text{Pfx}(W)$  dove  $W$  denota la stringa più lunga in  $M$ , mentre se  $M$  è infinito  $\text{Pfx}(M) = \{1\}^*$  .

**C10d.04** Sono molto più articolate delle precedenti le relazioni  $u \preceq_p w$  ,  $u \preceq_s w$  ,  $u \preceq_i w$  ,  $u \preceq_u w$  per un dato alfabeto  $A$  di due o più caratteri.

Innanzitutto si osserva che esse sono relazioni di ordine parziale. Si osserva poi che tra le relazioni precedenti considerate come insiemi di coppie di stringhe su  $A$ , cioè come sottoinsiemi di  $A^* \times A^*$ , sussistono le seguenti inclusioni:

$$(1) \quad \emptyset \subset \preceq_p \cap \preceq_s \subset \preceq_p , \preceq_s \subset \preceq_p \cup \preceq_s \subset \preceq_i \subset \preceq_u .$$

Insieme a queste relazioni interessano le loro associate riduzioni nonriflessive (univocamente determinate) che sono relazioni d'ordine stretto; le denoteremo, risp., con

$$\prec_p , \prec_s , \prec_i , \prec_u ;$$

Per esse useremo i termini, risp., prefisso proprio (o fattore sinistro proprio), suffisso proprio (o fattore destro proprio), infisso proprio o fattore proprio e sottostringa propria o sottosequenza propria.

Per esse vale una relazione di inclusione analoga alla (1):

$$(2) \quad \emptyset \subset \prec_p \cap \prec_s \subset \prec_p , \prec_s \subset \prec_p \cup \prec_s \subset \prec_i \subset \prec_u .$$

Inoltre sono spesso utili le corrispondenti relazioni di predecessore immediato, ovvero le loro riduzioni radice-star, e le rispettive relazioni trasposte di successore immediato.

Limitiamoci a considerare solo le relazioni di predecessore immediato; le denoteremo, risp., con  $\prec_{pI}$  ,  $\prec_{sI}$  ,  $\prec_{iI}$  e  $\prec_{uI}$  e per esse useremo i termini, risp., **prefisso immediato**, **suffisso immediato**, **infisso immediato** e **sottostringa immediata**.

Anche per esse vale una relazione di inclusione analoga alla (1):

$$\emptyset \subset \prec_{pI} \cap \prec_{sI} \subset \prec_{pI} , \prec_{sI} \subset \prec_{pI} \cup \prec_{sI} = \prec_{iI} \subset \prec_{uI} .$$

Queste relazioni si constata che sono fornite dalle seguenti espressioni:

$$\prec_{pI} = \bigcup_{w \in A^*} w \times (w, A) \quad , \quad \prec_{sI} = \bigcup_{w \in A^*} w \times (A, w) \quad ,$$

$$\prec_{\mathbf{pI}} \cup \prec_{\mathbf{sI}} = \mu \times \mathbf{A} \dot{\cup} \bigcup_{w \in \mathbf{A}^*} (w \times (w, \mathbf{A})) \cup (w \times (\mathbf{A}, w)) \quad ,$$

$$\prec_{\mathbf{pI}} \cap \prec_{\mathbf{sI}} = \{ \langle a, n \rangle \in \mathbf{A} \times \mathbb{N} \mid \langle a^n, a^{n+1} \rangle \} \quad , \quad \prec_{\mathbf{iI}} = \bigcup_{w \in \mathbf{A}^*} w \times \{w, \mathbf{A} \cup \mathbf{A}, w\} \quad .$$

**C10d.05** Per ogni  $w \in \mathbf{A}^+$  gli insiemi  $w \prec_{\mathbf{pI}}$  e  $w \prec_{\mathbf{sI}}$  sono costituiti da una sola stringa: quindi  $\mathbf{A}^*$  munito, risp., dell'insieme di archi  $\dot{\bigcup}_{w \in \mathbf{A}^*} \left( \bigcup_{x \in w} \succ_{\mathbf{pI}} \langle w, x \rangle \right)$  e dell'insieme di archi  $\dot{\bigcup}_{w \in \mathbf{A}^*} \left( \bigcup_{x \in w} \succ_{\mathbf{sI}} \langle w, x \rangle \right)$  costituisce due arborescenze illimitate; esse di solito vengono raffigurate svilupparsi verso l'alto [D30a].

//input pC10d04

//input pC10d04B

Munendo il terreno  $\mathbf{A}^*$  con l'insieme  $\dot{\bigcup}_{w \in \mathbf{A}^*} \left( \bigcup_{x \in w} \succ_{\mathbf{iI}} \langle w, x \rangle \right)$  si ottiene un digrafo.N graduato che non è un'arborescenza, in quanto per ogni  $w \in \mathbf{A}^+$  l'insieme  $\prec_{\mathbf{iI}}(w)$  contiene una sola stringa sse  $|w^{\text{minalf}}| = 1$ , mentre in caso contrario ne contiene almeno  $|w^{\text{minalf}}|$ .

//input pC10d05C

//input pC10d05D

Il terreno  $\mathbf{A}^+$  munito degli archi corrispondenti alle coppie costituenti l'ordine  $\succ_{\mathbf{uI}}$  conduce a un digrafo.N graduato che si può ottenere dal digrafo.N che corrisponde a  $\succ_{\mathbf{iI}}$  per aggiunta di numerosi archi.

//input pC10d05E

**(1) Eserc.** Relativamente a  $\mathbf{A} = \{a, b, c\}$ , tracciare i primi quattro livelli, cioè fino a raggiungere le stringhe di lunghezza 3, dei digrafi.N rappresentanti  $\succ_{\mathbf{iI}}$  e  $\succ_{\mathbf{uI}}$ .

**C10d.06** Della operazione di riflessione si può considerare l'estensione cartesiana applicata alle coppie o alle sequenze di stringhe e successivamente l'estensione booleana di questa estensione cartesiana. Utilizzando lo stesso simbolo per queste estensioni per le trasformazioni delle sequenze di stringhe e linguaggi abbiamo

$$\langle w, x \rangle^{\leftarrow} := \langle w^{\leftarrow}, x^{\leftarrow} \rangle \quad , \quad \langle \mathbf{L}, \mathbf{M} \rangle^{\leftarrow} := \langle \mathbf{L}^{\leftarrow}, \mathbf{M}^{\leftarrow} \rangle .$$

$$\langle w_1, w_2, \dots, w_n \rangle^{\leftarrow} := \langle w_1^{\leftarrow}, w_2^{\leftarrow}, \dots, w_n^{\leftarrow} \rangle \quad , \quad \langle \mathbf{L}_1, \mathbf{L}_2, \dots, \mathbf{L}_n \rangle^{\leftarrow} := \langle \mathbf{L}_1^{\leftarrow}, \mathbf{L}_2^{\leftarrow}, \dots, \mathbf{L}_n^{\leftarrow} \rangle .$$

Per l'estensione della riflessione alle coppie di stringhe e di linguaggi sarà utilizzato preferibilmente il termine specifico **biriflessione**.

Si verifica facilmente che le relazioni d'ordine  $\preceq_{\mathbf{p}}$  e  $\preceq_{\mathbf{s}}$  sono l'una la biriflessa dell'altra e che  $\preceq_{\mathbf{i}}$  e  $\preceq_{\mathbf{u}}$  sono invarianti per biriflessione.

Questo discende dalle implicazioni sulle coppie di stringhe:

$$w \preceq_{\mathbf{p}} x \iff (w^{\leftarrow}) \preceq_{\mathbf{s}} (x^{\leftarrow}) \quad , \quad w \preceq_{\mathbf{i}} x \iff (w^{\leftarrow}) \preceq_{\mathbf{i}} (x^{\leftarrow}) \quad , \quad w \preceq_{\mathbf{u}} x \iff (w^{\leftarrow}) \preceq_{\mathbf{u}} (x^{\leftarrow}) .$$

Può servire anche l'operazione unaria sulle coppie di stringhe o di linguaggi consistente nello scambio dei due membri della coppia stessa. Per questa operazione usiamo il termine trasposizione e scriviamo

$$\langle w, x \rangle^{\top} := \langle x, w \rangle \quad , \quad \langle \mathbf{L}, \mathbf{M} \rangle^{\top} = \langle \mathbf{M}, \mathbf{L} \rangle .$$

Queste due operazioni unarie sono evidentemente delle involuzioni e possono essere applicate successivamente ottenendo

$$(\langle w, x \rangle^{\leftarrow})^{\top} = \left( \langle w, x \rangle^{\top} \right)^{\leftarrow} = \langle x^{\leftarrow}, w^{\leftarrow} \rangle \quad , \quad (\langle \mathbf{L}, \mathbf{M} \rangle^{\leftarrow})^{\top} = \left( \langle \mathbf{L}, \mathbf{M} \rangle^{\top} \right)^{\leftarrow} = \langle \mathbf{M}^{\leftarrow}, \mathbf{L}^{\leftarrow} \rangle .$$

**C10d.07** In varie applicazioni risulta utile affiancare a ciascuno dei 4 ordini parziali di  $\mathbf{A}^*$  visti in precedenza l'ordine opposto e gli ordini che differiscono da questi due solo in quanto ottenuti procedendo sulle stringhe da destra a sinistra invece che da sinistra a destra.

Molti algoritmi riguardanti stringhe e linguaggi per essere ben definiti richiedono che si faccia riferimento a un ordine totale delle stringhe che estenda uno degli ordini parziali precedenti  $\preceq_{\mathbf{p}}$ ,  $\preceq_{\mathbf{s}}$ ,  $\preceq_{\mathbf{i}}$  e  $\preceq_{\mathbf{u}}$ .

Se denotiamo uno degli ordini precedenti con  $\preceq_{\mathbf{y}}$  e il richiesto ordine totale con  $\sqsubseteq$ , chiediamo che sia

$$w \preceq_{\mathbf{y}} x \implies w \sqsubseteq x ,$$

mentre non chiediamo il viceversa:  $w \sqsubseteq x \not\implies w \preceq_{\mathbf{y}} x$ .

Segnaliamo che un teorema di Garret Birkhoff garantisce che risulta possibile l'estensione di ogni ordine parziale sopra un terreno finito o numerabile a ordine totale compatibile.

Introduciamo due ordini totali basati sopra un ordine totale  $\leq$  dell'alfabeto  $\mathbf{A}$ , cioè basati su una sequenzializzazione delle lettere di  $\mathbf{A}$  che enunciamo scrivendo  $\langle \mathbf{a}_1 < \mathbf{a}_2 < \dots < \mathbf{a}_n \rangle$ .

Cominciamo con il confronto tra due stringhe della stessa lunghezza  $w =_{\mathbf{A}} \mathbf{a}_{j_1} \dots \mathbf{a}_{j_r}$  e  $x =_{\mathbf{A}} \mathbf{a}_{h_1} \dots \mathbf{a}_{h_r}$ ; si dice che  $w$  **precede secondo l'ordine lessicografico**  $x$ , e si scrive  $w <_{l_{xg}} x$ , sse si trova  $k \in \{1, \dots, r\}$  tale che  $i < k \implies \mathbf{a}_{j_i} = \mathbf{a}_{h_i}$  e  $\mathbf{a}_{j_k} < \mathbf{a}_{h_k}$ .

Per esempio, riferendoci all'usuale ordinamento alfabetico  $\langle \mathbf{a} < \mathbf{b} < \mathbf{c} < \dots < \mathbf{z} \rangle$ , abbiamo

$$\mathbf{a c b a z} <_{l_{xg}} \mathbf{a c b b a} <_{l_{xg}} \mathbf{a c c a a} .$$

**C10d.08** Facciamo riferimento a un alfabeto  $\mathbf{A}$  ordinato dalla richiesta  $\langle \mathbf{a}_1 < \mathbf{a}_2 < \dots < \mathbf{a}_n \rangle$  e ricordiamo che per ogni  $r \in \mathbb{P}$  le stringhe di  $\mathbf{A}^r$  sono ordinate totalmente dalla relazione  $\leq_{l_{xg}}$ .

Risulta semplice passare da una stringa di questo insieme a quella immediatamente precedente o a quella immediatamente successiva e conseguentemente risulta possibile procedere alla visita di  $\mathbf{A}^r$  secondo l'ordinamento  $\leq_{l_{xg}}$  o secondo l'ordine totale opposto  $\geq_{l_{xg}}$ .

Descriviamo l'algoritmo della suddetta visita di  $\mathbf{A}^r$  denotando con  $w = \mathbf{a}_{i_1} \mathbf{a}_{i_2} \dots \mathbf{a}_{i_r}$  la stringa di  $\mathbf{A}^r$  attualmente visitata.

- [1] si genera come prima stringa  $\mathbf{a}_1^r$ ;
- [2] se la lettera nella posizione più a destra della  $w$   $\mathbf{a}_{i_r} =: \mathbf{a}_s$  non è  $\mathbf{a}_n$ , l'ultima in  $\mathbf{A}$ , la si sostituisce con la successiva,  $\mathbf{a}_{s+1}$ ; quindi si emette come nuova stringa quella ottenuta e si torna alla [2];
- [3] se  $\mathbf{a}_{i_r} = \mathbf{a}_n$  si arretra cercando la posizione più a destra della  $w$ ,  $j$ , nella quale si abbia  $\mathbf{a}_{i_j} =: \mathbf{a}_t < \mathbf{a}_n$ ;
- [4] se si trova questo  $j \geq 1$ , si pone  $\mathbf{a}_{t+1}$  nella posizione  $j$  e si pone  $\mathbf{a}_1$  nelle posizioni successive  $h = j, j+1, \dots, r$ ; quindi si emette come nuova stringa quella ottenuta e si torna alla [2];

[5] se non si trova  $j$  la  $w$  precedentemente ottenuta è l'ultima di  $A^r$  e si arresta l'esecuzione.

L'ordine totale lessicografico si estende all'intero  $A^*$  mediante il seguente algoritmo di confronto di due stringhe generiche  $w = a_{j_1} \dots a_{j_r}$  e  $x = a_{h_1} \dots a_{h_s}$ :

- [1] scorrendo le due stringhe in parallelo si cerca il primo  $k \leq \min(r, s)$  tale che  $i < k \implies a_{j_i} = a_{h_i}$  e  $a_{j_k} \neq a_{h_k}$ ;
- [2] se non si trova tale  $k$  sono possibili tre situazioni:
  - . se  $w = x$  si decide  $w =_{lxg} x$ ;
  - . se  $w$  è prefisso proprio di  $x$  si decide  $w <_{lxg} x$ ;
  - . se  $x$  è prefisso di  $w$  si decide  $w >_{lxg} x$ ;
- [3] quando si incontra il suddetto  $k$ :
  - . se  $a_{j_k} < a_{h_k}$  si decide  $w <_{lxg} x$ ,
  - . se  $a_{j_k} > a_{h_k}$  si stabilisce  $w >_{lxg} x$ .

**C10d.09 Eserc.** Consideriamo  $u, v \in A^+$ ; si constati che

- (a)  $u <_{lxg} v \iff v \in uA^+ \vee u =: w \mathbf{a} x, v =: w \mathbf{b} y$  con  $w, x, y \in A^*, \mathbf{a}, \mathbf{b} \in A, \mathbf{a} < \mathbf{b}$ .
- (b)  $\forall w \in A^* : u <_{lxg} v \iff uw <_{lxg} vw$ .
- (c)  $u <_{lxg} v \implies \lceil \forall x, y \in A^* : ux <_{lxg} vy \rceil$ .

L'ordinamento lessicografico è un ordine totale che estende l'ordinamento parziale  $\preceq_p$ ; il suo primo elemento è  $\mu$  e la stringa immediatamente successiva secondo  $\leq_{lxg}$  di una generica  $w$  è  $w \mathbf{a}_1$ .

Non è invece possibile individuare una stringa che precede immediatamente secondo  $<_{lxg}$  una qualsiasi parola di  $A^* \setminus \{\mu, \mathbf{a}_1\}$ : infatti tra due stringhe  $u \mathbf{a}_i x$  e  $u \mathbf{a}_h y$  con  $\mathbf{a}_i < \mathbf{a}_h$  si trovano le infinite stringhe della forma  $u \mathbf{a}_i \mathbf{a}_j A^*$  per tutti gli  $\mathbf{a}_j < \mathbf{a}_h$  mentre tra  $w$  e  $w \mathbf{a}_1^m \mathbf{a}_k y$  con  $k \geq 1, \mathbf{a}_k > \mathbf{a}_1$  e  $y \in A^*$  si trovano le infinite stringhe dell'insieme  $w \mathbf{a}_1^m \mathbf{a}_1 A^*$ .

Quindi l'ordinamento lessicografico non è un ordine sequenziale.

**C10d.10** Un ordinamento sequenziale di  $A^*$  che estende  $\preceq_p$  e che risulta spesso utile è l'ordinamento **l'ordinamento secondo lunghezza-lessicografico** che denotiamo con  $\leq_{llx}$ .

Esso si basa sul seguente algoritmo di confronto di due generiche stringhe diverse  $w$  e  $x$ .

- [1] Primariamente si confrontano le lunghezze di  $w$  e  $x$ ; se  $w^{\lrcorner} < x^{\lrcorner}$  si decide  $w <_{llx} x$ , mentre se  $w^{\lrcorner} > x^{\lrcorner}$  si stabilisce  $w >_{llx} x$ ;
- [2] se le due lunghezze coincidono si effettua il loro confronto lessicografico e si decide di conseguenza, cioè si pone  $w <_{llx} x$  sse  $w <_{lxg} x$ .

L'ordine  $\leq_{llx}$  è un ordine sequenziale ed è agevole passare da una stringa  $w$  a quella immediatamente precedente o successiva. ¶ Per passare dalla  $w$  alla successiva, se  $w$  è l'ultima secondo  $\leq_{lxg}$  tra quelle aventi la sua stessa lunghezza  $r := w^{\lrcorner}$  si passa alla prima della lunghezza successiva, cioè alla  $\mathbf{a}_1^{r+1}$ ; se invece non è l'ultima si passa alla immediatamente successiva secondo  $\leq_{lxg}$  tra quelle della stessa lunghezza.

Per passare dalla  $w$  diversa dalla stringa muta alla precedente, sse  $w$  è la prima secondo  $\leq_{lxg}$  tra quelle aventi la sua stessa lunghezza  $r := w^{\lrcorner}$  si passa all'ultima della lunghezza precedente, cioè alla  $\mathbf{a}_r^{r-1}$ ; se invece non è la prima si passa alla immediatamente precedente secondo  $\leq_{lxg}$  tra quelle della stessa lunghezza  $r$ .

Evidentemente risulta semplice anche organizzare una procedura che a partire da una qualsiasi  $w \in A^*$  avanzi (o retroceda) secondo  $\leq_{llx}$  di un qualsiasi numero  $t$  di passi.

**C10d.11** Vediamo in dettaglio una procedura di visita illimitata che pone in corrispondenza biunivoca  $A^*$  con l'insieme degli interi naturali  $\mathbb{N}$ .

Per questo introduciamo la notazione  $\text{sopi}_{n,r}$  con  $r \in \mathbb{N}$  per esprimere la somma delle potenze di  $n$   $n^i$  per  $i = 0, 1, 2, \dots, r$  e osserviamo che tale somma esprime il cardinale delle stringhe su  $A$  con  $|A| = n$  aventi lunghezza minore o uguale a  $r$ , ossia

$$\text{sopi}_{n,r} := 1 + n + n^2 + \dots + n^r = |A^{\leq r}|.$$

Inoltre utilizziamo la variabile intera  $m$  da utilizzare come fosse una variabile del linguaggio  $c++$  con il ruolo di indice corrente delle stringhe successivamente visitate.

- [1] Si inizia con  $\mu$ , cui si fa corrispondere l'intero  $m = 0$ ;
- [2] si procede primariamente sulle successive lunghezze delle stringhe  $r = 1, 2, 3, \dots$ ;
  - [2.1] nell'ambito di un  $A^r$  si comincia con  $\mathbf{a}_1^r$  e le si associa  $\text{sopi}_{n,r-1}$ , cioè l'intero corrente  $m$  precedentemente aumentato di 1 (in  $c++$  si scriverebbe  $++m$ );
  - [2.2] si passa da una generica  $u \mathbf{a}_i \mathbf{a}_n^q$  con  $i = 1, 2, \dots, n-1$  e  $q = 0, \dots, |u| - 1$  alla immediatamente successiva  $u \mathbf{a}_{i+1} \mathbf{a}_1^q$ ; a questa si associa  $++m$ ;
  - [2.3] quando si è di fronte ad  $\mathbf{a}_n^r$ , ultima stringa di  $A^r$ , si aumenta la lunghezza passando ad  $\mathbf{a}_1^{r+1}$ , prima stringa di  $A^{r+1}$ , le si associa  $++m$  e si passa a [2.2].

La corrispondenza biunivoca che si stabilisce in tal modo, oltre ad associare 0 alla  $\mu$ , alla stringa  $\mathbf{a}_{j_1} \dots \mathbf{a}_{j_r}$  associa

$$(1) \quad m = \text{sopi}_{n,r-1} + j_1 n^{r-1} + j_2 n^{r-2} + j_3 n^{r-3} + \dots + j_{r-1} n + j_r.$$

La stringa di  $A^*$  che in tal modo si associa biunivocamente a un numero naturale  $m$  si dice che lo rappresenta, fedelmente, nel **sistema di numerazione  $n$ -adico**.

Si osserva che se  $A = \{1\}$  questa formula alla stringa " $|^k$ " associa l'intero naturale  $k$ , in accordo con la rappresentazione unadica degli interi naturali.

**C10d.12** Definiamo **derivata da sinistra rispetto a una stringa  $w$**  o **erosione da sinistra** di  $w$  dal linguaggio  $L$  il linguaggio

$$w \parallel L := \{z \in A^* \mid w \cdot z \in L\}.$$

Definiamo anche la **funzione o piccolo del linguaggio  $L$**  come

$$\mathbf{o}(L) := \begin{cases} \mu & \text{sse } \mu \in L, \\ \emptyset & \text{altrimenti.} \end{cases}$$

Per la precedente operazione e per la precedente funzione si trovano le proprietà che seguono.

- (1) **Prop.:**  $\forall v, w \in A^* : (v \cdot w) \parallel L = w \parallel (v \parallel L)$  ■
- (2) **Prop.:**  $\forall \mathbf{a}_{j_1} \dots \mathbf{a}_{j_k} \in A^+ : (\mathbf{a}_{j_1} \mathbf{a}_{j_2} \dots \mathbf{a}_{j_k}) \parallel L = \mathbf{a}_{j_k} \parallel (\dots(\mathbf{a}_{j_2} \parallel (\mathbf{a}_{j_1} \parallel L))\dots)$  ■
- (3) **Prop.:**  $w \in L \iff \mathbf{o}(w \parallel L) = \mu$  ■

**C10d.13** Si definisce **derivata da destra rispetto a una stringa  $z$**  o **erosione da destra** di  $z$  dal linguaggio  $L$

$$L \parallel z := \{w \in A^* \mid w \cdot z \in L\}.$$

Questa operazione è la duale-LR della derivazione da sinistra; infatti si ha

$$\forall z \in A^*, L \subseteq A^* : L \parallel z = (z^{\leftarrow} \parallel L^{\leftarrow})^{\leftarrow}.$$

Valgono quindi le uguaglianze duali per riflessione di quelle mostrate in d12.

- (1) **Prop.:**  $\forall v, w \in A^* : L \parallel (v \cdot w) = (L \parallel w) \parallel v$  ■

(2) Prop.:  $\forall \mathbf{a}_{j_1}, \dots, \mathbf{a}_{j_k} \in \mathbf{A} : L // (\mathbf{a}_{j_1} \mathbf{a}_{j_2} \dots \mathbf{a}_{j_k}) = ((\dots(L // \mathbf{a}_{j_k}) // \dots // \mathbf{a}_{j_2}) // \mathbf{a}_{j_1}) \blacksquare$

(3) Prop.:  $w \in L \iff \mathbf{o}(L // w) = \mu \blacksquare$

**C10d.14** Comprensibilmente risultano utili anche le estensioni booleane delle derivazioni rispetto a stringhe.

Si dice **derivata da sinistra rispetto a un linguaggio X** di L :

$$(1) \quad X \parallel L := \bigcup_{x \in X} x \parallel L .$$

Si dice **derivata da destra rispetto a un linguaggio Y** di L :

$$(2) \quad L // Y := \bigcup_{y \in Y} L // y .$$

(3) Prop.:  $X \parallel L = \{w \in \mathbf{A}^* \mid Xw \cap L \neq \emptyset\}$  ;  $L // Y = \{w \in \mathbf{A}^* \mid wY \cap L \neq \emptyset\} \blacksquare$

(4) Prop.:  $L // Y = ((Y^\leftarrow) \parallel (L^\leftarrow))^\leftarrow$  ,  $X \parallel L = ((L^\leftarrow) // (X^\leftarrow))^\leftarrow \blacksquare$

(5) Prop.:  $\mathbf{o}(X \parallel L) = \mu \iff \mathbf{o}(L // X) = \mu \iff X \cap L \neq \emptyset \iff \mathbf{o}(L // X) = \mu \iff \mathbf{o}(X // L) = \mu \blacksquare$

(6) Prop.:  $X \parallel L \neq \emptyset \iff X\mathbf{A}^* \cap L \neq \emptyset$  ,  $L // Y \neq \emptyset \iff L \cap \mathbf{A}^*Y \neq \emptyset \blacksquare$

(7) Prop.:  $X \parallel LM = (X \parallel L)M + (L // X) \parallel M$  ,  $(LM) // Y = L(M // Y) + L // (Y // M) \blacksquare$

**C10d.15** Con le operazioni di derivazione si possono esprimere prefissi, suffissi e infissi.

Per chiarire questi collegamenti è comodo servirsi delle scritte abbreviate

$$w // s := w // \mathbf{A}^s \quad \text{e} \quad s \parallel w := \mathbf{A}^s \parallel w \quad , \quad \text{per } s \in \mathbb{N} .$$

Inoltre consideriamo  $r \in \mathbb{P}$  e  $w =_{\mathbf{A}} \mathbf{a}_{j_1} \dots \mathbf{a}_{j_r}$  .

(1) Prop.:

$$\forall s \in [r] : w // s = \mathbf{a}_{j_1} \dots \mathbf{a}_{j_{r-s}} \quad , \quad s \parallel w = \mathbf{a}_{j_{s+1}} \dots \mathbf{a}_{j_r} \quad , \quad w = (w // s) \cdot ((r-s) \parallel w) \blacksquare$$

(2) Prop.:

$$\forall s \in [r] , t \in [r-s] : (s \parallel w) // t = s \parallel (w // t) = \mathbf{a}_{j_{s+1}} \dots \mathbf{a}_{j_{r-t}} \blacksquare$$

(3) Prop.:

$$\forall s \in [r] : s \parallel w = \mu \iff w // s = \mu \iff \text{len}(w) = s \blacksquare$$

(4) Prop.:

$$\forall s \in [r] : s \parallel w = \emptyset \iff w // s = \emptyset \iff \text{len}(w) < s \blacksquare$$

Possono essere utili anche le estensioni booleane delle abbreviazioni precedenti

$$s \parallel L := \mathbf{A}^s \parallel L \quad \text{e} \quad L // s := L // \mathbf{A}^s .$$

(5) Prop.: Consideriamo una stringa  $w$  e il relativo alfabeto minimo  $\mathbf{A}$ .

L'insieme dei prefissi di  $w$  è  $\mathbf{Pfx}(w) = w // \mathbf{A}^*$  ;

l'insieme dei suffissi di  $w$  è  $\mathbf{Sfx}(w) = \mathbf{A}^* \parallel w$  ;

l'insieme degli infissi di  $w$  è  $\mathbf{Ifx}(w) = \mathbf{A}^* \parallel (w // \mathbf{A}^*) = (\mathbf{A}^* \parallel w) // \mathbf{A}^* \blacksquare$

(6) Prop.: Consideriamo un linguaggio  $L$  e il relativo alfabeto minimo  $\mathbf{A}$ .

L'insieme dei prefissi di  $L$  è  $L // \mathbf{A}^*$  ;

l'insieme dei suffissi di  $L$  è  $\mathbf{A}^* \parallel L$  ;

l'insieme degli infissi di  $L$  è  $\mathbf{A}^* \parallel (L // \mathbf{A}^*) = (\mathbf{A}^* \parallel L) // \mathbf{A}^* \blacksquare$

L'ultima uguaglianza consente di introdurre la semplificazione

$$A^* \parallel L \parallel A^* := A^* \parallel L \parallel A^* = (A^* \parallel L \parallel A^* .$$

**C10d.16 (1) Eserc.** Mostrare che il numero degli elementi di  $\succ_{\mathbf{uI} \text{ } \mathbf{R}} w$  soddisfa le relazioni

$$|A| + 1 \leq |\succ_{\mathbf{uI} \text{ } \mathbf{R}} w| \leq |A| \cdot (w^{\dagger} + 1) .$$

**(2) Eserc.** Mostrare che, per  $L$  e  $M$  linguaggi qualsiasi:

$$\mathbf{o}(L^+) = \mathbf{o}(L) \quad , \quad \mathbf{o}(L^*) = \mu \quad , \quad \mathbf{o}(L + M) = \mathbf{o}(L) + \mathbf{o}(M) \quad , \quad \mathbf{o}(L \cdot M) = \mathbf{o}(L) \cdot \mathbf{o}(M) .$$

**C10d.17 (1) Eserc.** Dimostrare le seguenti proprietà della derivata da sinistra rispetto a linguaggi:

$$X \parallel (LM) = (X \parallel L) \cdot M + L \cdot ((L \parallel X) \parallel M) \quad , \quad (XY) \parallel L = Y \parallel (X \parallel L) \quad ,$$

$$X \parallel (L + M) = X \parallel L + X \parallel M \quad , \quad (X + Y) \parallel L = X \parallel L + Y \parallel L \quad ,$$

$$\mu \parallel L = L \quad , \quad X \parallel \mu = \mathbf{o}(X) \quad ,$$

$$X \parallel L^+ = ((L^* \parallel X) \parallel L) \cdot L^* \quad , \quad X \parallel L^* = \mathbf{o}(X) + ((L^* \parallel X) \parallel L) \cdot L^* \quad .$$

**C10d.18 (2) Eserc.** (a) Trovare le relazioni corrispondenti alle precedenti per la derivazione da destra.

(b) Constatare che  $\forall X, L, Y \in \mathbf{Lng} : X \parallel (L \parallel Y) = (X \parallel L) \parallel Y$  e conseguentemente giustificare l'abbreviazione  $X \parallel L \parallel Y := X \parallel (L \parallel Y)$ .

(c) Assunti  $L := a^2 b^2 + (ab)^3 + (ab^2)^2$ ,  $X := ab^* a$  e  $Y := aa^* b$ , trovare una significativa espressione per il linguaggio  $X \parallel L \parallel Y$ .

## C10 e. fattorizzazioni di stringhe

**C10e.01** Una stringa  $w \in A^+$  si dice **stringa primitiva** o anche **parola primitiva**, sse non si può esprimere come potenza di un'altra stringa (ovviamente più corta), cioè sse non si trova alcuna  $z \in A^+$  per la quale sia  $w = z^k$  per  $k \in \{2, 3, 4, \dots\}$ .

Una stringa primitiva e una nonprimitiva sono, risp.:

$$abbcabccabbacacbac \quad \text{e} \quad acbacbacbacbacb = (acb)^5.$$

Come mostra questo esempio, a parità di lunghezza, le stringhe nonprimitive si possono giudicare più regolari delle primitive della stessa lunghezza.

Questa maggiore regolarità si può chiamare anche minore complessità.

Più operativamente si stabilisce che una stringa  $w$  va considerata poco complessa quando si trova una espressione  $E$  che consente di individuarla per la quale sia  $|E| < |w|$ .

Questo significa che le stringhe nonprimitive possono essere individuate fornendo meno informazioni di quante richiedono le primitive della stessa lunghezza. In altre parole le stringhe primitive sono meno abbreviabili delle nonprimitive.

Denotiamo con  $\text{Primw}_A$  il linguaggio delle stringhe primitive su  $A$ . Per esempio

$$\begin{aligned} \text{Primw}_{\{a,b\}} = \{ & a, b, ab, ba, aab, aba, abb, baa, bab, bba, \\ & aaab, aaba, aabb, abaa, abba, abbb, baaa, baab, babb, bbab, bbaa, bbab, bbba, \dots \} \end{aligned}$$

Ovviamente per ogni alfabeto  $A$  l'insieme di stringhe  $\text{Primw}_A$  è invariante per permutazione dei caratteri di  $A$ .

**C10e.02** Vediamo alcuni fatti collegati alla non primitività delle stringhe.

**(1) Prop.:** Le stringhe di lunghezza 1 sono tutte primitive, ovviamente.

Le stringhe la cui lunghezza è un numero primo aut sono primitive aut sono potenze di una lettera ■

**(2) Prop.:** Per ogni  $w \in A^+$  si può effettivamente individuare una sola parola primitiva  $z$  tale che  $w \in z^+$ .

**Dim.:** Se  $w^{\#} = 1$   $w \in A$ . Se  $w^{\#}$  è un numero primo si vede se è una potenza della sua prima lettera e in caso contrario è primitiva.

Consideriamo quindi una  $w$  la cui lunghezza  $s := \text{len}(w)$  abbia sottomultipli propri maggiori di 1.

Organizziamo l'esame dei successivi prefissi di  $w$  aventi come lunghezza uno dei suddetti sottomultipli procedendo per sottomultipli crescenti.

Per ciascuno di questi prefissi, che denotiamo con  $z$ , si controlla se soddisfa il requisito  $z^{w^{\#}/z^{\#}} = w$ .

Il primo prefisso  $z$  che soddisfa questo requisito è la stringa primitiva cercata: infatti ogni ulteriore prefisso  $y$  tale che  $y^{|w|/|y|} = w$  deve essere una potenza di  $z$ ; infatti

$$y^{|w|/|y|} = z^{|w|/|z|} \iff y^{|w|/|z|} = z^{|w|/|y|} \iff y = z^{|y|/|z|}.$$

Se non si trova alcun prefisso che soddisfa il requisito la  $w$  è primitiva ■

**C10e.03 Prop.** Due stringhe diverse commutano sse sono potenze di una stessa stringa (primitiva).

**Dim.:** Consideriamo due stringhe  $w, x \in A^+$ .

Se  $w = z^n$  e  $x = z^m$  per qualche stringa  $z$  è ovvio che sia  $wx = z^{m+n} = xw$ .

Quindi basta dimostrare  $w x = x w \implies A^+ \ni z \sqsupseteq w, x \in z^+$ , cosa che faremo costruttivamente con l'algoritmo che segue che permette di precisare come ottenere la  $z$  a partire da  $w$  e  $x$ .

**Algoritmo:** Non può essere  $w^{\perp} = x^{\perp}$ , in quanto la loro commutazione implica la loro uguaglianza, esclusa per ipotesi.

In caso contrario basta esaminare il caso  $w^{\perp} < x^{\perp}$ . È allora possibile scrivere  $x = w^h v$  con  $h \in \mathbb{P}$  e  $0 \leq v^{\perp}$ . Se  $v^{\perp} = 0$  come stringa richiesta si può scegliere  $z := w$  e  $x = z^h$ . Se  $0 < v^{\perp}$  si ha  $w^{h+1} v = w^h v w$  e quindi  $w v = v w$  con  $v^{\perp} < x^{\perp}$ .

Siamo quindi ricondotti a una ricerca del tipo di quella iniziale, ma limitata a coppie di stringhe aventi somma delle lunghezze inferiore.

Si può quindi procedere con successive fasi che ripetono le manovre precedenti con la  $w$  rimpiazzata da una stringa  $v$  via via più corta e dopo un numero finito di fasi si individua la  $z$  tale che  $w, v \in z^+$  e quindi anche tale che  $x = w^h v \in z^+$  ■

**C10e.04 Prop.** Consideriamo due stringhe  $w, x \in A^+$ .

Si trova una coppia di interi positivi  $\langle m, n \rangle$  tali che sia  $w^m = x^n \iff A^+ \ni z$  tale che  $w, x \in z^+$ .

**Dim.:** “ $\Leftarrow$ ”: è implicazione ovvia.

“ $\Rightarrow$ ”: Se  $m = n$  si assume  $z := w = x$ . Se  $m = 1$  si assume  $z := x$  e, dualmente-LR, se  $n = 1$  si assume  $z := w$ .

Se  $m$  ed  $n$  hanno un divisore comune  $k > 1$ , si può scrivere  $m = kp$  ed  $n = kq$ ; in tal caso  $w^m = x^n$  equivale a  $w^p = x^q$ .

Resta quindi da dimostrare  $\lceil w^m = x^n \implies A^+ \ni z$  tale che  $w, x \in z^+$  per  $m$  ed  $n$  mutuamente primi e, data l'equivalenza dei ruoli di  $w$  e  $x$ , possiamo limitarci a supporre  $m < n$ , ovvero  $w^{\perp} > x^{\perp}$ .

Si può scrivere  $w = x^h y$  con  $0 < y^{\perp} < x^{\perp}$  e  $h \in \mathbb{P}$ ; dall'esame dell'inizio di  $w$  si ottiene  $y x^n = x^n y$  e quindi, dato che  $y^{\perp} < x^{\perp}$ , dopo eventuali successive riduzioni dell'uguaglianza precedente, si ottiene  $x y = y x$ ; per questa e utilizzando l'algoritmo in e03 per questa uguaglianza si può individuare la stringa  $z$  tale che  $x, y \in z^+$  e  $w = x^h y \in z^+$  ■

**C10e.05** Riprendiamo la nozione di occorrenza di carattere in una stringa, nozione che aiuta a precisare i rapporti tra una stringa e i suoi infissi, le sue eventuali regolarità e permette di definire due utili relazioni di equivalenza nei monoidi liberi.

Assumiamo ancora  $A = \{a_1, \dots, a_n\}$  e Consideriamo la stringa  $w =_A w(1) \dots w(s) \in A^s$ .

Ogni coppia  $\langle a_j, i \rangle \in A \times \{s\}$  si dice **occorrenza di un carattere**  $a_j$  nella stringa  $w$  sse  $w(i) = a_j$ .

Più in generale ogni coppia  $\langle u, i \rangle \in A^{\leq s} \times \{s\}$  si dice **occorrenza dell'infisso**  $u$  nella  $w$  sse  $w = x u y$  con  $x^{\perp} = i - 1$ .

Denotiamo con  $\text{Occ}(u, w)$  l'insieme delle occorrenze di  $u$  nella  $w$  e denotiamo con  $|w|_u$  o con  $\text{occ}(u, w)$  il numero delle occorrenze di  $u$  nella stringa  $w$ .

**(1) Prop.:**  $\sum_{a_i \in A} |w|_{a_i} = s = |w|$  ■

**(2) Prop.:**  $|w|_u > 0 \iff \text{Occ}(u, w) \neq \emptyset \iff u \in \text{Ifx}(w)$  ■

Studiando le parole e i linguaggi su  $A$  può essere significativo anche l'alfabeto minimo di una parola  $w$ , i.e. l'insieme delle lettere che occorrono in  $w$ , ossia l'insieme delle  $a_i \in A \sqsupseteq |w|_{a_i} > 0$ ; tale alfabeto lo denotiamo con  $w^{\text{minalf}}$ .

Palesemente  $w^{\text{minalf}} = \text{Ifx}(w) \cap A$ .

**C10e.06** Assegnamo ad  $A$  un ordine totale  $<$ , cioè poniamo in sequenza i caratteri che costituiscono questo alfabeto; più precisamente supponiamo che sia  $\mathbf{a}_1 < \mathbf{a}_2 < \dots < \mathbf{a}_n$ .

Diciamo **vettore di Parikh della stringa**  $w$  relativo ad  $A$  e  $\leq$  la sequenza di lunghezza  $n = |A|$  che segue

$$\mathit{Prk}(w) = w^{\mathit{Prk}} := \langle |w|_{\mathbf{a}_1}, |w|_{\mathbf{a}_2}, \dots, |w|_{\mathbf{a}_n} \rangle .$$

Per esempio relativamente all'alfabeto  $\{a, b, c, d\}$  ordinato nel modo usuale risulta:

$$(\mathit{accdadaca})^{\mathit{Prk}} = \langle 4, 0, 3, 2 \rangle .$$

La **funzione di Parikh** “ $\mathit{Prk}$ ” associa a ogni stringa di un  $A^*$  un elemento di  $\mathbb{N}^{|A|}$ , è suriettiva e noniniettiva sse  $|A| > 1$ .

In particolare è evidente che per ogni stringa (anche nonpalindroma)  $w$  si ha  $\mathit{Prk}(w^\leftarrow) = \mathit{Prk}(w)$ .

La relativa equivalenza, che definiamo  $\sim_{\mathit{Prk}} := \mathit{Prk} \circ_{lr} (\mathit{Prk}^{-1})$ , è costituita dalle coppie di stringhe di  $A^*$  ottenibili l'una dall'altra permutando le occorrenze di carattere che le compongono.

Per esempio  $(\mathit{accdadaca})^{\mathit{Prk}} = (a^4 c^2 d c d)^{\mathit{Prk}} = (c^2 a d c a^2 d a)^{\mathit{Prk}} = \langle 4, 0, 3, 2 \rangle$  e quindi  $\mathit{accdadaca} \sim_{\mathit{Prk}} a^4 c^2 d c d \sim_{\mathit{Prk}} c^2 a d c a^2 d a$ .

Queste classi corrispondono ai multiinsiemi basati sull'insieme ordinato  $A$ ; esse quindi sono in biiezione con permutazioni con ripetizioni di  $A$  [B13e17] e i loro cardinali sono dati da coefficienti multinomiali. Per esempio il numero di stringhe equivalenti ad “*essere*” è  $6!/(3!2!1!) = 720/(6 \cdot 2) = 60$ , mentre la sequenza in ordine lessicografico delle stringhe dell'insieme  $\mathit{essere}_A \sim_{\mathit{Prk}}$  comprende:

*eeerss eeesrs eersss eeress eerses eersse ereess ereses eresse ersees ersese erssee*  
 .....  
 .....  
 .....  
*seseer sesere sesree sreees sreese sresee srseee ssreee sseree sseere sseree ssreee* .

**C10e.07** Le classi dell'equivalenza  $\sim_{\mathit{Prk}}$  sono interamente contenute negli insiemi  $A^s$  di stringhe aventi la stessa lunghezza; in altri termini la partizione di  $A^*$  relativa a  $\sim_{\mathit{Prk}}$  è un raffinamento di quella relativa alla uguaglianza delle lunghezze.

In genere è utile individuare le classi di una equivalenza attraverso loro elementi rappresentativi facilmente individuabili. Un atteggiamento di questo tipo spesso conveniente consiste nell'assumere come rappresentativo di ciascuna classe l'elemento minimo secondo un opportuno ordine totale dell'insieme ambiente.

Come elementi rappresentativi delle classi di  $\sim_{\mathit{Prk}}$ , come per i rappresentativi di altre equivalenze entro insiemi  $A^* = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}^*$ , conviene assumere le stringhe minime secondo l'ordine lessicografico  $\leq_{lsg}$ : la generica classe  $w_A \sim_{\mathit{Prk}}$  viene dunque rappresentata da

$$\mathbf{a}_1^{|w|_{\mathbf{a}_1}} \mathbf{a}_2^{|w|_{\mathbf{a}_2}} \dots \mathbf{a}_n^{|w|_{\mathbf{a}_n}} .$$

Dal punto di vista dell'insieme dei vettori di Parikh  $\mathbb{N}^{\times n}$  si osserva che

$$\forall w, x \in A^* : (w, x)^{\mathit{Prk}} = w^{\mathit{Prk}} \mathbf{+} x^{\mathit{Prk}} ,$$

dove con  $\mathbf{+}$  denotiamo la somma termine a termine delle  $n$ -uple di  $\mathbb{N}^{\times n}$ ; la funzione di Parikh quindi è un epimorfismo di monoidi da  $\langle A^*, \cdot, \mu \rangle$  su  $\langle \mathbb{N}^n, \mathbf{+}, \mathbf{0}_n \rangle$ .

Anche della funzione di Parikh risulta utile la estensione booleana; essa è un morfismo tra algebre di Kleene e c'è da aspettarsi che gli insiemi di sequenze di interi naturali ottenibili dalla applicazione di

questa funzione del genere  $[\text{Lng}_A \mapsto \mathbb{N}^{\times |A|}]$  ad alcuni linguaggi forniscano caratteristiche interessanti dei linguaggi stessi.

**C10e.08** Una proprietà che contribuisce alla regolarità delle stringhe, più debole dell'essere esprimibile come potenza di una stringa più corta [e01, ..., e04] riguarda la possibilità di trovare in esse occorrenze ripetute.

Si dice **coppia sovrapposta** sull'alfabeto  $A$  una stringa su  $A$  che può scriversi  $w = v y = x v$  per qualche  $v, x, y \in A^+$ .

Un esempio di coppia sovrapposta è **ABRACADABRA**; per essa  $v = \text{ABRA}$ .

Si dice invece **stringa con overlap** su  $A$  una stringa della forma  $v u v u v$  con  $v \in A^*$  ed  $u \in A^+$ ; l'infisso ripetuto  $v u v$  si dice **overlap della stringa**  $v u v u v$ .

Una stringa con overlap fornisce un caso particolare di coppia sovrapposta: infatti si può scrivere  $v u v u v = w y = z w$  con  $w = v u v$ ,  $y = u v$  e  $z = v u$ .

Si osserva che le stringhe che sono quadrati di stringa, cioè che hanno la forma  $u u$ , sono particolari stringhe con overlap ( $v = \mu$ ); tali sono anche tutte le stringhe potenze  $k$ -esime di stringhe con  $k = 3, 4, \dots$ , cioè stringhe della forma  $u^k$ ; per una tale stringa vi sono  $k - 1$  scelte dell'overlap:  $u, u^2, \dots, u^{k-1}$ .

**C10e.09 (1) Prop.:** Ogni coppia sovrapposta  $w = v y = x v$  sull'alfabeto  $A$  presenta un overlap come prefisso e un overlap come suffisso.

**Dim.:** Possiamo limitarci al caso in cui sia  $x^{\perp} < v^{\perp}$ , dopo aver osservato che il caso  $x^{\perp} = v^{\perp}$  si riduce a  $x = v$  e dopo aver segnalato la possibilità di trattare il caso  $x^{\perp} > v^{\perp}$  per dualità-LR.

Poniamo  $z := x \parallel v$ , in modo che sia  $w = x z y = x v = v y$  e osserviamo che  $x^{\perp} < v^{\perp}$  implica  $|z| > 0$ .

Denotiamo con  $a$  il primo carattere di  $v$ , osserviamo che  $a$  è anche il carattere iniziale di  $z$ ; quindi introduciamo  $z_1 := a \parallel z$  e osserviamo che  $z_1$  è suffisso di  $z, v$  e  $w$ .

Posto  $t := v \parallel z_1$  si ha  $w = a t a t a z_1$  e quindi si trova una stringa con overlap che è prefisso di  $t$ .

La dimostrazione della presenza di un overlap suffisso si ottiene per dualità-LR, ossia procedendo come sopra, ma servendosi delle stringhe riflesse di quelle usate in precedenza ■

**C10e.10** Consideriamo una parola  $w \in A^s$  che contiene due occorrenze di  $v \in A^r$  con  $0 < 2r \leq s$  che scriviamo  $\langle v, i \rangle$  e  $\langle v, i' \rangle$ ; scriviamo inoltre  $w = x v y = x' v y'$  con  $x^{\perp} = i - 1$  e  $x'^{\perp} = i' - 1$ .

Si danno tre possibilità per le occorrenze di  $v$ :

- si hanno **occorrenze sconnesse** sse  $i + r < i'$ , cioè sse si può scrivere  $w = x v z v y'$  con  $z \in A^+$ ;
- si hanno **occorrenze adiacenti** sse  $i + r = i'$ , cioè sse si può scrivere  $w = x v v y'$ ;
- si hanno **occorrenze sovrapposte** sse  $i + r > i'$ , cioè sse si può scrivere  $w = x v' y'$  con  $v'$  coppia sovrapposta.

Si può stabilire che la presenza di occorrenze ripetute in una parola contribuisca alla sua regolarità e che questa sia accentuata se nella parola sono presenti occorrenze adiacenti, cioè fattori esprimibili come quadrati o come potenze superiori a due.

**(1) Prop.:** Una stringa  $w$  ha come infisso una coppia sovrapposta  $\iff w$  ha come infisso un overlap.

**Dim.:** " $\implies$ ": Dalla e04 segue che la presenza di una coppia sovrapposta implica la presenza di un overlap .

" $\impliedby$ ": Come si è detto in e08(1), la presenza di un overlap in una parola costituisce un caso particolare di presenza di una coppia sovrapposta ■

**C10e.11** Due stringhe  $w$  ed  $x$  si dicono **stringhe ciclicamente coniugate** sse una di esse si può ottenere permutando ciclicamente l'altra; per segnalare questa situazione scriviamo  $w \sim_{\mathcal{C}} x$ .

Chiaramente  $w \sim_{\mathcal{C}} x$  sse  $w$  ed  $x$  si possono fattorizzare, risp., come  $w = uv$  e  $x = vu$  sse in  $A^+$  si trova una stringa  $u$  tale che  $wu = ux$ .

La partizione di  $A^*$  criptomorfa all'equivalenza  $\sim_{\mathcal{C}}$  è un raffinamento della partizione criptomorfa all'equivalenza  $\sim_{Prk}$ : infatti  $\sim_{Prk}$  aggrega stringhe ottenibili attraverso permutazioni qualsiasi dei caratteri componenti, mentre  $\sim_{\mathcal{C}}$  aggrega a una stringa solo quelle ottenibili attraverso permutazioni cicliche. Per esempio  $abc \sim_{Prk} acb$ , mentre  $abc \not\sim_{\mathcal{C}} acb$ .

Constatiamo che le classi di  $\sim_{\mathcal{C}}$  contenute in  $\{a, b\}^3$  sono  $\{a^3\}$ ,  $\{b^3\}$ ,  $\{aab, aba, baa\}$  e  $\{abb, bab, bba\}$  e che esse coincidono con le classi di  $\sim_{Prk}$ .

Le classi dell'equivalenza  $\sim_{\mathcal{C}}$  contenute in  $\{a, b\}^4$  sono  $\{a^4\}$ ,  $\{b^4\}$ ,  $\{aaab, aaba, abaa, baaa\}$ ,  $\{aabb, abba, bbaa, baab\}$ ,  $\{abbb, babb, bbab, bbba\}$  e  $\{abab, baba\}$ ; ora il loro insieme non coincide con l'insieme delle classi relative a  $\sim_{Prk}$  dal quale si ottengono sostituendo la quarta e la sesta classe con la loro unione  $\{aabb, abba, bbaa, baab, abab, baba\}$ .

Sia le classi di  $\sim_{\mathcal{C}}$  che le classi di  $\sim_{Prk}$  sono opportunamente rappresentate dalle rispettive stringhe minime secondo un ordine lessicografico.

Nell'ambito di  $\{a, b\}^4$  questo sistema di rappresentativi per  $\sim_{\mathcal{C}}$  è  $\{a^4, aaab, aabb, abab, aaab, b^4\}$ , mentre per  $\sim_{Prk}$  è  $\{a^4, aaab, aabb, abbb, b^4\}$ .

Può essere chiarificante scrivere

$$w \sim_{\mathcal{C}} x \implies Prk(w) = Prk(x), \quad \text{mentre} \quad Prk(w) = Prk(x) \not\implies w \sim_{\mathcal{C}} x.$$

**C10e.12** Si osserva che se una stringa  $w$  è potenza di una stringa più corta, cioè se  $w = z^q$  con  $q > 1$ , ogni stringa ottenuta permutando ciclicamente  $w$  è potenza  $q$ -esima di una permutazione ciclica di  $z$ . Quindi tra le classi di  $\sim_{\mathcal{C}}$  si distinguono quelle costituite solo da stringhe primitive dalle classi formate da stringhe che sono tutte potenze proprie di stringhe più corte.

Volendo esprimere le varie stringhe di  $A^*$  come prodotti di stringhe più corte, le stringhe primitive risultano più interessanti delle nonprimitive in quanto consentono maggiore concisione.

Vedremo nella sezione che segue che tra le primitive risultano particolarmente utili le minime lessicografiche delle classi di  $\sim_{\mathcal{C}}$ .

## C10 f. parole di Lyndon

**C10f.01** Si dice **parola di Lyndon** sull'alfabeto  $A$  munito di un ordine totale  $<$  una stringa di  $A^+$  che è primitiva ed è la minima secondo l'ordine lessicografico nella propria classe di  $\sim_C$ , classe di stringhe ciclicamente coniugate, ossia inferiore- $l_{xg}$  alle sue rotazioni.

Denotiamo con  $\text{Lynd}_{A,<}$  il linguaggio delle parole di Lyndon per  $\langle A, < \rangle$ . Quando  $<$  si può considerare implicito, per esempio se  $A$  è costituito dalle prime lettere di un alfabeto ben riconosciuto, usiamo la scrittura semplificata  $\text{Lynd}_A$ .

In particolare possiamo scrivere:

$$\text{Lynd}_{\{a,b\}} = \{a, b, ab, aab, abb, aaab, aabb, abbb, \\ aaaaab, aaabb, aabab, aabbb, ababb, abbbb, aaaaab, \dots\}$$

$$\text{Lynd}_{\{a,b,c\}} = \{a, b, c, ab, ac, bc, aab, aac, abb, abc, acc, bbc, bcc, \\ aaab, aaac, aabb, aabc, aacb, aacc, abac, abbb, abbc, abcb, abcc, acbb, acbc, acce, \\ aaaaab, aaabb, aaabc, aaacb, aaacc, aabab, aabac, aabbb, aabbc, aabcb, aabcc, \\ aacab, aacac, aacb, aacbc, aaccb, aaccc, ababb, ababc, abacc, abbab, abbac, \\ abbbb, abbbc, abccb, abccc, abeac, abebb, abebc, abccb, abccc, acacb, acacc, \\ acbbb, acbbe, acbcc, acbbb, accbc, acccb, acccc, \\ bbbbc, bbccc, bbcbc, bbccc, bcbbc, beccc, aaaaab, \dots\}.$$

**C10f.02** Sono evidenti gli enunciati che seguono riguardanti alfabeti  $A$  e  $B$ .

- (1) **Prop.:** (a)  $\text{Lynd}_1 = \{1\}$ .  
 (b)  $A \subset B \implies \text{Lynd}_A \subset \text{Lynd}_B$ .  
 (c)  $|\text{Lynd}_A| = \aleph_0$ .  
 (d)  $\text{Lynd}_{\{a,b\}} \supset \{b\} \dot{\cup} a^+ b^+$ .  
 (e)  $A = \{a_1 < a_2 < \dots < a_r\} \implies a_1^+ a_2^+ \dots a_r^+ \subset \text{Lynd}_A$ .  
 (f)  $\text{Lynd}_A = A \dot{\cup} \text{Lynd}_A \cap AA^+$  ■

**C10f.03 (1) Prop.:** Una stringa  $w \in AA^+ = A^{\geq 2}$  appartiene a  $\text{Lynd}_A$  sse precede strettamente secondo l'ordine lessicografico ogni suo suffisso proprio, ovvero

$$\text{Lynd}_A = \left\{ w \in A^+ \mid \forall x \in w^{\text{Sfx}} \cap A^+ : w \leq_{l_{xg}} x \right\}.$$

**Dim.:** Introduciamo l'insieme delle stringhe su  $A$  che precedono secondo  $<_{l_{xg}}$  ogni loro suffisso proprio, ovvero poniamo

$$\mathbf{L} := \left\{ w \in A^+ \mid \forall x \in w^{\text{Sfx}} \cap A^+ : w \leq_{l_{xg}} x \right\}.$$

Osserviamo che la proprietà vale per le lettere di  $A \subset \text{Lynd}_A$  in quanto non posseggono suffissi propri; per le stringhe di lunghezza maggiore o uguale a 2 procediamo a dimostrare due relazioni di inclusioni opposte.

$\lceil \text{Lynd}_A \subseteq \mathbf{L} \rceil$ : Consideriamo  $w \in \mathbf{L}$ , la stringa  $x \in w^{\text{Sfx}} \cap A^+$  tale che sia  $w <_{l_{xg}} x$  ed introduciamo  $u := w // x$ , ossia  $u \in A^+$  tale che  $w = ux$ .

Procediamo a dimostrare, per assurdo, che  $x \notin w^{\text{Pfx}}$ . Se fosse  $w = xt$  con  $t \in A^+$  sarebbe  $xt = ux$  e quindi, grazie a e08, si trovano  $y, z \in A^*$  ed  $i \in \mathbb{N}$  tali che si avrebbe

$$u = yz, \quad t = zy, \quad (yz)^i y \quad \text{e quindi} \quad w = (yz)^{i+1} y.$$

La parola  $w \in \text{Primw}$  sarebbe inferiore per  $<_{l_{xg}}$  di ogni sua coniugata ciclica, cioè sarebbe

$$(yz)^{i+1}y <_{lxg} y(yz)^{i+1} \quad \text{e quindi} \quad , (zy)^{i+1} <_{lxg} (yz)^{i+1} .$$

Giustapponendo la  $y$  a destra ai due membri si avrebbe  $(zy)^{i+1}y <_{lxg} (yz)^{i+1}y = w$ , contraddicendo la precedenza lessicografica di ciascuna delle parole di Lyndon rispetto alle sue coniugate cicliche. Dunque ogni iparola di Lyndon  $w$  appartiene ad  $\mathbf{L}$ .

$\lceil \mathbf{L} \subseteq \text{Lynd}_{\mathbf{A}} \rceil$  : Consideriamo  $w \in \mathbf{L} \cap \mathbf{A}\mathbf{A}^+$  e le stringhe  $u, v \in \mathbf{A}^+$  tali che  $uv = w$  e si abbia la disuguaglianza  $w <_{lxg} v$ ; ovviamente  $w <_{lxg} vu$  e quindi  $w \in \text{Lynd}_{\mathbf{A}}$  ■

**C10f.04** Dato che  $\mathbf{A} \subset \text{Lynd}_{\mathbf{A}}$ , evidentemente ogni stringa su  $\mathbf{A}$  si può esprimere come prodotto di parole di Lyndon. Le fattorizzazioni che si servono di fattori di Lyndon sono dette **fattorizzazioni di Lyndon**; esse godono di proprietà notevoli.

Cominciamo con l'osservare esplicitamente che una parola di Lyndon può presentare diverse fattorizzazioni di Lyndon. Due esempi:

$$aabb = aab,b = a,abb \quad , \quad ababb = ab,abb = a,b,abb$$

**(1) Prop.:**  $l, m \in \text{Lynd}_{\mathbf{A}} \wedge l <_{lxg} m \implies lm \in \text{Lynd}$  .

**Dim.:** Chiaramente  $lm <_{lxg} m$ ; se  $l <_{\text{pfx}} m$  scrivendo  $m = lm'$  abbiamo  $m <_{lxg} m'$  e di conseguenza  $lm \leq_{lxg} lm' = m$ ; se invece  $l \notin m\mathbf{A}^*$ , essendo  $l <_{\text{pfx}} m$ , si ha ancora  $lm \leq_{lxg} lm' = m$ .

Prendiamo ora una  $v \in m^{\text{Sfx}}$ : dato che  $m \in \text{Lynd}_{\mathbf{A}}$ , grazie a f03(1), abbiamo  $lm \leq_{lxg} m \leq_{lxg} v$ . Introduciamo anche  $v' \in l^{\text{Sfx}}$ : si ha  $l \leq_{lxg} v'$  e quindi  $lm \leq_{lxg} v'm$ .

Dunque  $lm$  precede lessicograficamente ogni suo suffisso proprio e, grazie a f03(1)  $lm \in \text{Lynd}_{\mathbf{A}}$  ■

**(2) Prop.:** Sia  $w \in \text{Lynd}_{\mathbf{A}} \setminus \mathbf{A}$  e sia  $m$  la più lunga delle parole in  $w^{\text{Sfx}} \cap \text{Lynd}_{\mathbf{A}} \setminus \{w\}$ , insieme che contiene almeno l'ultima lettera di  $w$ . Introdotto  $l$  tale che  $w = lm$ , allora  $l \in \text{Lynd}_{\mathbf{A}}$  e  $l <_{lxg} lm <_{lxg} m$  .

**Dim.:** Se  $l \in \mathbf{A}$  l'enunciato vale. In caso contrario scegliamo  $v \in l^{\text{Sfx}}$ . Dato che  $vm \notin \text{Lynd}_{\mathbf{A}}$ , scegliamo  $t \in (vm)^{\text{Sfx}} \setminus \{vm\}$  tale che  $t <_{lxg} vm$ .

Allora se  $v <_{lxg} t$ , essendo  $v <_{lxg} t <_{lxg} vm$  si deduce che esiste  $s <_{lxg} m$  tale che  $t = vs$  e questa  $s$  sarebbe un suffisso proprio di  $m$  lessicograficamente inferiore ad  $m$ , contro la  $m \in \text{Lynd}_{\mathbf{A}}$ ; dunque  $t \leq_{lxg} v$ ; di conseguenza  $l <_{lxg} lm <_{lxg} t \leq_{lxg} v$ , consentendo di affermare  $l <_{lxg} v$  e, grazie a f03(1),  $l \in \text{Lynd}_{\mathbf{A}}$ .

Se invece  $l <_{lxg} m$ , dato che  $w = lm \in \text{Lynd}_{\mathbf{A}}$ , si ottiene  $lm <_{lxg} m$  ■

Gli enunciati (1) e (2) conducono ad un'altra caratterizzazione delle parole di Lyndon.

**(3) Prop.:**  $w \in \text{Lynd}_{\mathbf{A}} \setminus \mathbf{A} \iff w = lm \quad \text{con} \quad l, m \in \text{Lynd}_{\mathbf{A}} \wedge l <_{lxg} m$  ■

**C10f.05** L'enunciato f04(1) suggerisce una procedura ricorsiva in grado di procedere illimitatamente nella generazione delle parole di un  $\text{Lynd}_{\mathbf{A}}$  a partire dalle lettere dell'alfabeto  $\mathbf{A}$ .

Questa procedura presenta una iterazione primaria (illimitata) che procede a generare le parole di Lyndon delle successive lunghezze  $\ell = 2, 3, 4, \dots$

In ogni stadio relativo ad una  $\ell$  si ha una iterazione secondaria che riguarda coppie di parole già acquisite che presentano coppie di lunghezze  $\langle h, k \rangle$  con  $h = 1, \dots, \ell - 1$  e  $k = \ell - h$ .

Una iterazione terziaria riguarda le coppie di parole  $\langle l, m \rangle$  relative ad ogni attuale coppia  $\langle h, k \rangle$  e per ciascuna di queste  $\langle l, m \rangle$  verifica se  $l <_{lxg} m$  ed in caso affermativo accoda  $lm$  alla successione delle parole di Lyndon.

Descriviamo le prime azioni della procedura nel caso  $\mathbf{A} = \{a, b\}$ . Per  $\ell = 2$  aggiunge alla lista delle lettere la sola  $a, b$ . Per  $\ell = 3$  emette  $a, ab$  ed  $ab, b$ . Per  $\ell = 4$  emette  $a, aab, a, abb, aab, b$  e  $abb, b$ .

.....

Si vede come si presentano anche parole ripetute e risulta chiara la convenienza di eliminare al più presto le ripetizioni per evitare successive manovre inutili.

La coppia  $\langle l, m \rangle \in \text{Lynd}_A^{\times 2}$  tale che  $w := lm \in \text{Lynd}_A$  ed  $m$  di lunghezza massimale tra i suffissi di  $w$  che sono parole di Lyndon si dice **fattorizzazione-mls** della  $w$  (o fattorizzazione standard, o fattorizzazione con suffisso di lunghezza massimale). Denotiamo con  $\text{fmls}(w)$  questa coppia (evidentemente unica).

**(1) Prop.:** Consideriamo  $w \in \text{Lynd}_A \cap AA^+$  e scriviamo  $\langle l, m \rangle := \text{fmls}(w)$ .

Allora  $\forall x \in \text{Lynd}_A \cap (w \text{ a } >_{lxg}) : \langle w, x \rangle = \text{fmls}(wx) \iff x <_{lxg} m \blacksquare$

### C10f.06 Teorema (teorema di Lyndon)

Ogni stringa  $w \in A^+$  possiede una unica fattorizzazione mediante una sequenza di parole di Lyndon lessicograficamente noncrescenti.

**Dim.:** Come già osservato [f04], ogni  $w \in A^+$  possiede almeno una fattorizzazione di Lyndon; scegliamone una la cui lunghezza  $n$  sia minimale e scriviamola  $w =: l_1 l_2 \dots l_n$ .

Se vi fosse  $i \in [n - 1]$  tale che  $l_i <_{lxg} l_{i+1}$ , grazie a f04(1) si avrebbe  $l_i l_{i+1} \in \text{Lynd}_A$  e quindi la fattorizzazione non sarebbe minimale. Quindi la generica  $w$  possiede almeno una fattorizzazione di Lyndon nondecrecente-lxg.

Dimostriamo per assurdo l'unicità di questa fattorizzazione supponendo che esista un'altra fattorizzazione di  $w$  con gli stessi requisiti e scriviamo

$$l_1 l_2 \dots l_n = l'_1 l'_2 \dots l'_n$$

$$\text{ove } \forall i \in [1, n] : l_i, l'_i \in \text{Lynd} \text{ e } \forall i \in [1, n - 1] : l_i \geq_{lxg} l_{i+1} \wedge l'_i \geq_{lxg} l'_{i+1}.$$

Se vi fosse un  $j \in [1, n - 1] \text{ s.t. } |l_j| > |l'_j|$ , si avrebbe  $l_j = l'_j \dots l'_{j+h} u$  con  $h \in \mathbb{N}$  e  $u \in l'_j$  **Sfx**.

Grazie a f03(1) si avrebbe allora  $l_j <_{lxg} u \leq_{lxg} l'_{j+h+1} \leq_{lxg} l'_j \leq_{lxg} l_j$ , un evidente assurdo.

Quindi  $l_i = l'_i \blacksquare$

**C10f.07** Anche la precedente dimostrazione suggerisce un algoritmo per fattorizzare una qualsiasi stringa mediante parole di Lyndon.

Un algoritmo più veloce è suggerito dal seguente fatto.

**(1) Prop.:** Sia  $l_1 l_2 \dots l_n$  una fattorizzazione della  $w \in A^+$  mediante una sequenza di parole di Lyndon noncrescente secondo  $<_{lxg}$ . Allora  $l_n$  è il suffisso di  $w$  minimo secondo  $<_{lxg}$ .

**Dim.:** Sia  $v$  il suffisso di  $w$  minimo-lxg e scriviamo  $w = uv$ ; in forza di f04(1) si ha  $v \in \text{Lynd}$ .

Se  $u = \mu$  l'asserto è dimostrato. Se no scriviamo  $s$  il suffisso di  $u$  minimo-lxg e scriviamo  $w =: rsv$ . Non può essere  $s <_{lxg} v$ , in quanto, grazie a f04(1), sarebbe  $sv \in \text{Lynd}$  da cui  $sv <_{lxg} v$ , contro la definizione di  $v$ ; dunque  $s \geq_{lxg} v$  e proseguendo questo processo si giunge alla enunciata fattorizzazione di  $w \blacksquare$

**C10f.08** Si osservi che da f07(1) segue direttamente f06.

Va anche segnalato che l'algoritmo suggerito da f07(1) non è il più veloce possibile.

Art M. Duval nel 1980 ha dimostrato che conviene procedere sulla  $w$  da fattorizzare non da destra verso sinistra, ma da sinistra a destra, in quanto in tal modo si ottiene la fattorizzazione con un numero di confronti tra lettere lineare nella lunghezza della stringa sottoposta.