

Grammar Induction

Apprendimento di strutture sintattiche tramite esempi:
applicazione al linguaggio naturale

Elaborazione del Linguaggio Naturale
A.A. 2003/2004

Daniela Lepri

Sommario

- Introduzione
- Grammar Induction: Teorie e Metodi
- Tecniche di apprendimento di Linguaggi Formali
- Tecniche di apprendimento del Linguaggio Naturale
- Risultati e Conclusioni
- Riferimenti

16/02/2007

2

Sommario

- **Introduzione**
 - Cos'è la Grammar Induction
 - Base induttiva: Stringhe e Sequenze
 - Paradigmi dei Linguaggi Formali
 - Paradigmi del Linguaggio Naturale
- Grammar Induction: Teorie e Metodi
- Tecniche di apprendimento di Linguaggi Formali
- Tecniche di apprendimento del Linguaggio Naturale
- Risultati e Conclusioni
- Riferimenti

16/02/2007

3

Cos'è la Grammar Induction

Una definizione informale:

Processo di apprendimento di grammatiche da un insieme di informazioni sul linguaggio.

- **Origini:**
Disciplina sviluppatasi intorno agli anni '60, a seguito degli studi di Gold [Gold, 1967] "Language Identification in the Limit".
- **Alias:**
 - "Grammatical Inference"
 - "Automatic Language Acquisition"

16/02/2007

4

Definizione Formale di Grammatica

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

N : Insieme finito dei simboli Non-Terminali

Σ : Insieme finito dei simboli Terminali

$S \in N$: Simbolo Non-Terminale iniziale

$R \subset (N \cup \Sigma)^* N (N \cup \Sigma)^* \times (N \cup \Sigma)^*$: Insieme di Regole

Notazione per le Regole di Produzione:

$$\alpha \rightarrow \beta, \quad \alpha \in (N \cup \Sigma)^* N (N \cup \Sigma)^*, \quad \beta \in (N \cup \Sigma)^*$$

$$\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots$$

16/02/2007

9

Grammatiche e Linguaggi

Derivazione Elementare:

$$\mu \alpha \delta \xrightarrow{G} \mu \beta \delta \quad \text{iff} \quad \exists (\alpha \rightarrow \beta) \in R, \quad \mu, \delta \in (N \cup \Sigma)^*$$

Derivazione $\xrightarrow{*}_G$: sequenza finita di derivazioni elementari

$D_G(x)$: insieme delle *Derivazioni* di una stringa $x \in \Sigma^*$ tale che $S \xrightarrow{*}_G x$

Una G è *ambigua* se $\exists x \in \Sigma^*$ tale che $|D_G(x)| > 1$

Linguaggio generato dalla Grammatica:

$$\mathcal{L}(G) = \{x \in \Sigma^* \mid S \xrightarrow{*}_G x\}$$

16/02/2007

10

Gerarchia di Chomsky

0: Unrestricted

1: Context Sensitive

$$\alpha \rightarrow \beta, \quad |\alpha| \leq |\beta|$$

➤ 2: Context Free (CFG)

$$B \rightarrow \beta, \quad B \in N$$

➤ 3: Regular, Finite-State (FS Automata) or Right Linear

$$A \rightarrow aB \text{ or } A \rightarrow a, \quad A, B \in N, \quad a \in \Sigma \cup \{\lambda\}$$

16/02/2007

11

Regular Grammars and Finite-State Automata

■ Regular Grammars (RG):

$$G = (N, \Sigma, R, S), \\ (A \rightarrow aB) \in R \vee (A \rightarrow a) \in R, \quad A, B \in N, \quad a \in \Sigma$$

■ Finite-State Automata (FSA):

$$\mathcal{A} = (Q, \Sigma, \delta, q_0, F), \quad q_0 \in Q, \quad F \subseteq Q, \quad \delta: Q \times \Sigma \rightarrow 2^Q$$

Risultato di equivalenza:

Per ogni RG esiste un FSA che riconosce lo stesso linguaggio

16/02/2007

12

Stochastic Grammars e Linguaggi

- Una *Stochastic Grammar* G' è la grammatica G , alla quale sono state associate probabilità ad ogni produzione

$$G' = (G, p), \quad G = (N, \Sigma, R, S), \quad p : R \rightarrow [0, 1]$$

- G' è *propria* se $\forall A \in N \quad \sum_{A \rightarrow \beta \in R} p(A \rightarrow \beta) = 1$

- Probabilità di una stringa x generata da G'

$$\forall x \in \Sigma^* \quad p(x|G') = \sum_{d \in D_{G'}(x)} p(d), \quad p(d) = \prod_{(A \rightarrow \beta) \in d} p(A \rightarrow \beta)$$

16/02/2007

13

Sommario

- Introduzione
 - Cos'è la Grammar Induction
 - Strutture Sintattiche: Esempi
 - Paradigmi dei Linguaggi Formali
 - **Paradigmi del Linguaggio Naturale**
 - Grammatiche che generano il Linguaggio Naturale
 - Perché Context Free
 - Esempio di CFG per il Linguaggio Naturale
- Grammar Induction: Teorie e Metodi
- Tecniche di apprendimento di Linguaggi Formali
- Tecniche di apprendimento del Linguaggio Naturale
- Risultati e Conclusioni
- Riferimenti

16/02/2007

14

Grammatiche che generano il Linguaggio Naturale

- Secondo il cosiddetto approccio "generativista", il linguaggio naturale può essere generato da una Grammatica Formale.
- Approccio fondato da Noam Chomsky
- Motivato dal fatto che le persone sono in grado di generare frasi che sono sintatticamente corrette e totalmente nuove (cioè che non hanno mai sentito prima)
- La Grammatica Generativa racchiude in sé le regole sintattiche del linguaggio naturale.

16/02/2007

15

Grammatiche che generano il Linguaggio Naturale(2)

- *Simboli Terminali*: Categorie Lessicali o part-of-speech (Verb, Noun, Adjective, Adverb, Preposition, Determiner)
 - Per comodità, i rispettivi lessemi, sono a volte inclusi nella grammatica come simboli terminali.
- *Simboli Non-Terminali*: Categorie Sintattiche (Sentence, Noun Phrase, Verb Phrase, Prepositional Phrase, Relative clause)
- Solitamente si utilizzano Grammatiche Context Free

16/02/2007

16

Perché Context Free

- E' provato che gli Automi a Stati Finiti non sono abbastanza potenti per modellare tutte le frasi di un Linguaggio Naturale

Ad Esempio:

- Centre-embedded sentences (Frase "annidate")
 - Esempio semplice: (The book (the professor wrote) is good)
 - Esempio più complicato: (The book (the professor (the students (who are doing well) like) wrote) is good)
- E' simile al linguaggio non regolare: $a^nbc^n \mid n > 0$
- Le CFG riescono a modellizzare queste strutture ricorsive.
- Una CFG per $a^nbc^n \mid n > 0$
 $\langle S \rangle ::= a \langle S \rangle c \mid b$

16/02/2007

17

Esempio di CFG per il Linguaggio Naturale

Esempio di CFG elementare per l'Inglese:

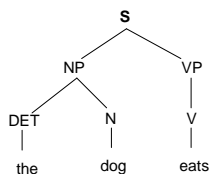
- | | |
|--|-------------|
| 1. S -> NP VP | Sentence |
| 2. NP -> DET N N | Noun Phrase |
| 3. VP -> V V NP | Verb Phrase |
| 4. DET -> the this that my many which... | Determiners |
| 5. N -> dog girl park... | Nouns |
| 6. V -> eats eat look... | Verbs |

16/02/2007

18

Esempio di generazione e Parse Tree: "the dog eats"

- Derivazione:
 $S \rightarrow NP VP \rightarrow DET N VP \rightarrow DET N V \rightarrow the\ dog\ eats$
- Tagging: [N, dog] [V, eats] [DET, the]
- Bracketing: [S [NP [DET, the] [N, dog]] [VP [V, eats]]]
- Parse Tree:



16/02/2007

19

Sommario

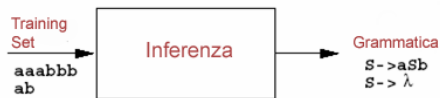
- Introduzione
- **Grammar Induction: Teorie e Metodi**
 - Definizione Formale di GI
 - "Language identification in the limit" e Computabilità
 - Classificazione dei metodi
- Tecniche di apprendimento di Linguaggi Formali
- Tecniche di apprendimento del Linguaggio Naturale
- Risultati e Conclusioni
- Riferimenti

16/02/2007

20

Definizione Formale di GI

- Teorie e metodi per apprendere induttivamente strutture sintattiche comuni ad un insieme di esempi e informazioni (*Training Data*), che si presume siano generati dallo stesso processo.



16/02/2007

21

Definizione Formale di GI (2)

- Le informazioni (*Training Data*) a disposizione sono solitamente campioni di stringhe definite sullo stesso alfabeto.
 - *Positive Set*: $S_+ \subset L$
 - *Negative Set*: $S_- \subset \Sigma^* - L$
- Problema: trovare G tale che:
 - $S_+ \subset \mathcal{L}(G)$
 - $S_- \cap \mathcal{L}(G) = \emptyset$

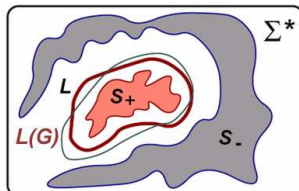
16/02/2007

22

Teorema di Gold: "Language Identification in the limit"

Condizione sufficiente per equivalenza tra il linguaggio generato dalla grammatica inferita $L(G)$ e il linguaggio target L . [Gold, 1967]

- Se il Training Set $S_+ \cup S_-$ è sufficientemente ampio e rappresentativo, allora si ha: $\mathcal{L}(G) = L$



" $L(G) \rightarrow L$ per $|S_+ \cup S_-| \rightarrow \infty$ "

16/02/2007

23

Computabilità del Teorema di Gold

- Ogni classe di linguaggi Recursively Enumerable (cioè, per i quali si può costruire una Turing Machine) può essere identificata al limite.
- Non sempre si conosce S_-
- Senza S_- non si identificano al limite nemmeno i linguaggi Regolari

16/02/2007

24

Classificazione dei Metodi

- Utilizzando solo S+
 - Problema: *sovra-generalizzazione*
 - Induzione di una grammatica il cui linguaggio è una generalizzazione (include strettamente) del linguaggio target.
 - Tale processo è irreversibile senza S-
- Utilizzando sia S+ che S-
 - Non si ha sovra-generalizzazione
 - I linguaggi Regolari sono Identificabili al Limite
 - I linguaggi Context-Free non lo sono
- Inferenza di Stochastic Grammars
 - Possono essere appresi efficacemente utilizzando unicamente S+
 - Controllano sovra-generalizzazione con informazioni statistiche dedotte dal Positive Set
 - Più usati in pratica

16/02/2007

25

Sommario

- Introduzione
- Grammar Induction: Teorie e Metodi
- **Tecniche di apprendimento di Linguaggi Formali**
 - Apprendimento di Finite-State Automata
 - Un Algoritmo Generico con S+ e S- e RPNI
 - Un Algoritmo Generico per indurre Stochastic Grammars
 - Applicazioni e Risultati
- Tecniche di apprendimento del Linguaggio Naturale
- Risultati e Conclusioni
- Riferimenti

16/02/2007

26

Apprendimento di Finite-State Automata

- Un primo contatto con gli approcci della GI
- Riconoscono Linguaggi Regolari
- Le proprietà dei Linguaggi Regolari sono ben note
- I Metodi di Inferenza sono più semplici da sviluppare e comprendere
- Ogni linguaggio può essere approssimato con precisione arbitraria con un Linguaggio Regolare (in senso stocastico)
- Molti problemi pratici sono modellati efficacemente da FSA

16/02/2007

27

Il Prefix Tree Acceptor

- Insieme dei prefissi di un linguaggio

$$L \subseteq \Sigma^* : Pr(L) = \{u \in \Sigma^* | uv \in L, v \in \Sigma^*\}$$

- Prefix Tree Acceptor sul Positive Set

$$PT(S_+) = (Q, S, \delta, q_0, F)$$

$$Q = Pr(S_+); q_0 = \lambda; F = S_+; \delta(u, a) = ua \text{ iff } u, ua \in Pr(S_+)$$

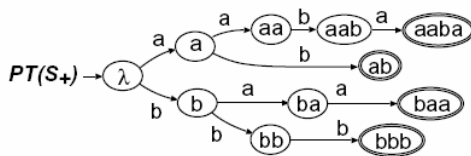
Automa che riconosce tutte e solo le stringhe contenute in S+

16/02/2007

28

PT(S+): Esempio

$S+ = \{ab, aaba, baa, bbb\}$



16/02/2007

29

Un Algoritmo Generico con S+ e S-

- 1. Costruisci il Prefix Tree Acceptor PT(S+), sul Positive Set.
- 2. Scegli due stati sui quali tentare una *merge* (fusione)
- 3. Effettua la merge su tali stati
- 4. Se esiste una stringa nel Negative Set, accettata dal nuovo automa, annulla la merge e ripristina l'automata precedente
- 5. Ripeti 2-3-4 fino a quando viene raggiunto il criterio di Stop (precisione desiderata)

16/02/2007

30

Regular Positive-Negative Inference (RPNI)

- Sviluppato da Oncina e Garcia [Oncina, 92]
- Utilizza sia S+ che S-
- Costruisce un Automa a Stati Finiti Deterministico, compatibile con S+ e S-, senza cercare di minimizzare l'automata stesso
- Riesce ad identificare efficientemente al limite, ogni linguaggio regolare.

16/02/2007

31

Algoritmo RPNI

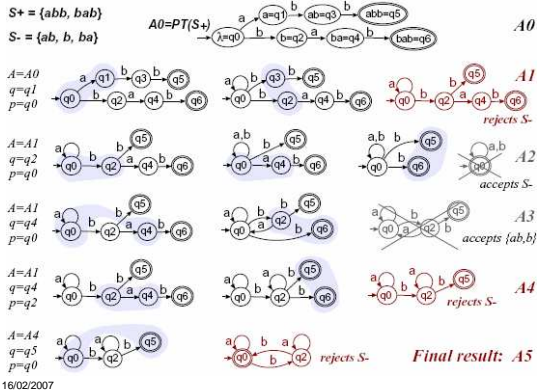
Algorithm RPNI (*Regular Positive & Negative Inference*)
Input: S+ S-
Output: A: DFA which accepts S+ and do not accept R-

Method: A=PT(S+); (let Q(A) denote the set of states of A)
 forall q in Q(A) in lexicographic order do
 forall p < q in lexicographic order do
 A'=merge(A,p,q)
 while A' is not deterministic do
 select q', q'' which violate determinism
 A'=merge(A',q',q'')
 endwhile
 if no string from S- is accepted by A'
 then A=A'
 endif
 end forall p
 end forall q

end RPNI
 16/02/2007

32

Algoritmo RPNI: Esempio

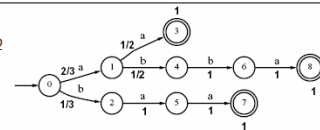


Un Algoritmo Generico per indurre Stochastic Grammars

Probabilistic Prefix Tree Acceptor (PPTA):

E' un PTA alle transizioni del quale sono associate delle probabilità

PPTA è *proprio*



1. Costruisci PPTA
2. **Seleziona** una coppia di stati (i, j)
3. Se (i, j) sono **Compatibili**, effettua la merge di i e j
4. Ripeti 2,3 fino a che l'Automa raggiunge il livello di precisione prestabilito

16/02/2007

34

Un Algoritmo Generico per Stochastic Grammars (2)

- Vari criteri di Selezione e Compatibilità possibili
 - Selezione della coppia di stati:
 - Ordine Breadth-first nel PPTA
 - Ordine Lessico-grafico degli stati
 - Compatibilità della coppia di stati:
 - Privilegiati stati con transizioni uscenti che hanno probabilità simili (simile quantità di stringhe riconosciute passando per quella transizione) (ALERGIA)
 - Privilegiate coppie che minimizzano il trade-off tra la perdita di precisione (probabilità di riconoscere il linguaggio target) e la diminuzione della dimensione dell'Automa (MFI)

16/02/2007

35

Applicazioni e Risultati

- Esempi di campi di applicazione con successo:
 - Speech Recognition
 - Riconoscimento della lingua parlata
 - Modellazione del Linguaggio
 - Elaborazione Digitale di Musica
 - Apprendimento e riconoscimento automatico di stili musicali
- Funzionano solo con Linguaggi Regolari
- Non applicabili per indurre Grammatiche del Linguaggio Naturale

16/02/2007

36

Sommario

- Introduzione
- Grammar Induction: Teorie e Metodi
- Tecniche di apprendimento di Linguaggi Formali
- **Tecniche di apprendimento del Linguaggio Naturale**
 - Classificazione dei metodi
 - Vari Approcci
- Risultati e Conclusioni
- Riferimenti

16/02/2007

37

Tecniche di Apprendimento del Linguaggio Naturale

- Inducono Grammatiche Context Free
- E' impossibile indurre CFG da un testo
 - conseguenza del risultato di Gold sul concetto di "Language identification in the limit"
- Alternativa: PAC Learning – Probably Approximately Correct Learning
 - Ci si accontenta di un'approssimazione che rispetti il criterio di precisione desiderato

16/02/2007

38

Tecniche di Apprendimento del Linguaggio Naturale (2)

- Vari approcci utilizzati:
 - Algoritmi genetici [Lankhorst, 1994]
 - Reti Neurali [Honkela, 1995]
 - Minimum Description Length Principle (MDL): minimizzano la Grammatica [Wolff, 1988]
 - Metodi Euristici basati sul merging di categorie e valutazione della probabilità di successo della Grammatica indotta: Algoritmo inside-outside [Baker, 1979], [Lari & Young, 1992]

16/02/2007

39

Classificazione dei metodi

- (Supervised) vs (Unsupervised)
 - I primi prevedono un intervento umano durante il processo di inferenza
 - - Generalmente operano su testi pre-tagged: richiedono lavoro aggiuntivo per annotare il testo
 - + Risultati migliori: più precisi grazie alle info aggiuntive
 - I secondi sono metodi totalmente automatizzati
 - + Generalmente operano senza informazioni aggiuntive sulla struttura del testo: meno costosi.
 - - Meno precisi: non hanno idea del linguaggio target
 - + Funzionano anche quando è impossibile annotare il testo (Es. dati non linguistici, DNA, codici cifrati)
- (Solo Positive Set) vs (Positive e Negative Set)
 - Negative Set solitamente costituito da frasi casuali

16/02/2007

40

Sommario

- Introduzione
- Grammar Induction: Teorie e Metodi
- Tecniche di apprendimento di Linguaggi Formali
- Tecniche di apprendimento del Linguaggio Naturale
 - Classificazione dei metodi
 - **Vari Approcci**
 - EMILE 4.1
 - An Evolutionary Approach (Algoritmo Genetico)
- Risultati e Conclusioni
- Riferimenti

16/02/2007

41

EMILE 4.1 [Adriaans, 2000]

- Basato sull'algoritmo di EMILE 3.0 sviluppato da Adriaans [Adriaans, 1992, 1999]
- Rientra nei paradigmi *PAC* (Probably Approximately Correct Learning)
- *Unsupervised*: uno dei primi algoritmi efficienti operanti su testo non annotato.
- S+: utilizza solo esempi positivi
- Costruisce una **Categorial Grammar**
- Assume che il Linguaggio Naturale sia **shallow** ("semplice")
 - Ogni costrutto sintattico ha una frase che lo esemplifica, e la quale lunghezza è di ordine logaritmico rispetto alla complessità della Grammatica inferita.

16/02/2007

42

Categorial Grammar

- Classe di grammatiche che generano tutti i linguaggi Context-Free
- Basate sull'assegnamento di *Categorie* sintattiche alle parole del lessico
 - Esempio: [Cat, NOUN]
- Una *Categoria* grammaticale è caratterizzata da:
 - *Espressioni* che appartengono alla categoria
 - *Contesti*, nei quali compaiono le suddette espressioni
 - Un contesto che appare in più *Categorie*, viene chiamato *Ambiguo*
- Esempio di coppia *contesto/espressione*:
 - "John (makes) tea"
 - Informalmente: L'*espressione* "makes" appare nel *contesto* "John (.) tea"
 - Regola di derivazione formale: makes \rightarrow John \ α / tea
 - " α A" = Espressione α può essere preceduta da A
 - " α B" = Espressione α può essere seguita da B
- Proprietà: *Principio di Sostituzione*
 - Espressioni della stessa categoria possono essere sostituite tra di loro in tutti i contesti appartenenti alla categoria stessa.
 - Esempio: "John (likes) tea"
- Utilizzate per comodità nell'ambito della Morfologia e della Linguistica Computazionale

16/02/2007

43

EMILE 4.1: Principi di funzionamento

L'induzione di EMILE si basa su due fatti:

- **Principio di sostituzione** delle Categorial Grammars
- In un testo sufficientemente vasto, ricorrono con frequenza le stesse combinazioni contesto/espressione (**cluster**)

16/02/2007

44

EMILE 4.1: Algoritmo

- 1. Estrae tutti i possibili *clusters* contesto/espressione dal Set di frasi
- 2. Raggruppa espressioni che appaiono nello stesso contesto
- 3. Raggruppa i contesti nei quali appaiono le stesse espressioni
- 4. Sceglie le categorie grammaticali in base al metodo *2-dimensional clustering*
- 5. Trasforma le categorie grammaticali così trovate in regole di derivazione

16/02/2007

45

EMILE: Esempio 1D-clustering

■ Frasi utilizzate

- "John makes tea", "John likes tea"

- 1. Estrae tutte le possibili combinazioni di coppie "context/expression"

	(.)	John	John	(.)	John	(.)	John	(.)
	makes	(.)	makes	(.)	makes	(.)	likes	likes
	tea	x	tea	(.)	tea	(.)	tea	(.)
John		x						x
makes			x					
tea				x				x
John makes					x			
makes tea						x		
John makes tea							x	
likes			x					
John likes					x			
likes tea						x		
John likes tea							x	

- 2. Raggruppa espressioni che compaiono nello stesso contesto

- Tabella context/expression per 1D-clustering

16/02/2007

46

EMILE: Esempio 1D-clustering (2)

■ Clusters contesto/espressione indotti:

- [{'makes', 'likes'}, 'John (.) tea']
- [{'John makes', 'John likes'}, '(.) tea']
- [{'makes tea', 'likes tea'}, 'John (.)']
- [{'John makes tea', 'John likes tea'}, '(.)']

16/02/2007

47

EMILE: Esempio 1D-clustering (3)

- 3. Raggruppa contesti, nei quali compaiono le stesse espressioni

- Per esemplificare il passo 3, consideriamo le frasi:

- "John makes tea"
- "John likes tea"
- "John makes coffee"
- "John likes coffee"

	John	John	John	John
	(.)	(.)	(.)	(.)
	tea	coffee	makes	likes
makes		x	x	
likes		x	x	
tea				x
coffee				x

- Parte rilevante della Tabella context/expression per 1D-clustering

- Clusters contesto/espressione indotti:

- [{'makes', 'likes'}, {'John (.) tea', 'John (.) coffee'}] ← Verb
- [{'tea', 'coffee'}, {'John makes (.)', 'John likes (.)'}] ← Noun

- Categorie Grammaticali (Verb, Noun) indotte direttamente dai Clusters

16/02/2007

48

EMILE: Problema 1D-Clustering

- Problema: 1D-Clustering non sa gestire contesti *ambigui*

Frase utilizzate

- "John makes tea", "John likes tea"
- "John makes coffee", "John likes coffee"
- "John is eating", "John likes eating"

	John (.) tea	John (.) coffee	John (.) eating	John makes (.)	John likes (.)	John is (.)
makes	x	x				
likes	x	x	x			
is			x			
tea				x	x	
coffee				x	x	
eating					x	x

- Tabella context/espression per 2D-clustering

16/02/2007

49

EMILE: Problema 1D-Clustering (2)

- Intuitivamente, si possono dedurre le seguenti Categorie Grammaticali:

- Noun-Phrases NP ('tea', 'coffee')
- Verb-Phrases VP ('makes', 'likes')
- 'Ing'-Phrases IP ('eating')
- Auxiliary Verbs AV ('is')

- Clusters contesto/espressione indotti con 1D-clustering:

```
[ {'John (.) tea', 'John (.) coffee'}, {'makes', 'likes'} ]
[ {'John (.) eating'}, {'likes', 'is'} ]
A → [ {'John makes (.)'}, {'tea', 'coffee'} ]
B → [ {'John likes (.)'}, {'tea', 'coffee', 'eating'} ]
[ {'John is (.)'}, {'eating'} ]
```

- NB: il contesto "John likes (.)" è ambiguo!

- Compare sia nella categoria Noun-Phrase che in quella Ing-Phrase

- I clusters "A" e "B" non possono essere raggruppati con 1D Clustering (poco naturale)

16/02/2007

50

EMILE: Esempio 2D-Clustering

- Matrice con blocchi di dimensione massima in evidenza

	John (.) tea	John (.) coffee	John (.) eating	John makes (.)	John likes (.)	John is (.)
makes	x	x				
likes	x	x	x			
is			x			
eating				x	x	
tea				x	x	
coffee				x	x	

- Per l'induzione dei clusters context/espression:

- Non è richiesto che il blocco sia totalmente contenuto nella matrice
- E' sufficiente che una certa percentuale di esso lo sia
 - E' improbabile che S+ contenga tutte le frasi necessarie a "riempire" un blocco

16/02/2007

51

EMILE: Esempio 2D-Clustering(2)

- Clusters contesto/espressione indotti:

```
[ {'John (.) tea', 'John (.) coffee'}, {'makes', 'likes'} ]
[ {'John (.) eating'}, {'likes', 'is'} ]
→ [ {'John makes (.)', 'John likes (.)'}, {'tea', 'coffee'} ]
[ {'John is (.)', 'John likes (.)'}, {'eating'} ]
[ {'John (.) tea', 'John (.) coffee', 'John (.) eating'}, {'likes'} ]
[ {'John likes (.)'}, {'eating', 'tea', 'coffee'} ]
```

- Induzione più vicina all'intuizione umana

- Gli ultimi 2 clusters sono ridondanti:

- Coppie context/espression già presenti nei cluster precedenti

- Alla fine del ciclo di esecuzione, EMILE rimuove questi clusters superflui.

16/02/2007

52

EMILE: Determinazione delle regole di derivazione

- Per ogni espressione e appartenente alla categoria T si ha: $[T] \rightarrow e$
- Sottoespressioni di un tipo che sono anche espressioni di un altro tipo, vengono sostituite con quest'ultimo, esempio:
 - EMILE ha già identificato la categoria
 - $[T1] \rightarrow \text{dog} \mid \text{cat} \mid \text{parrot}$
 - EMILE identifica le seguenti espressioni come appartenenti alla categoria T2
 - $\{\text{I feed my dog, I feed my cat, I feed my parrot}\}$
 - Regola inferita
 - $[T2] \rightarrow \text{I feed my } [T1]$

16/02/2007

53

EMILE: Esempio (5)

- Determinazione di regole di derivazione ricorsive
 - Le seguenti espressioni appartengono alla categoria [S]
 - "Mary drinks tea"
 - "John observes that Mary drinks tea"
 - Regola inferita
 - $[S] \rightarrow \text{John observes that } [S]$
 - Frase indotta
 - "John observes that John observes that Mary drinks tea"

16/02/2007

54

EMILE: The Medline experiments

- Esecuzione di EMILE sul dominio bio-medico, su un corpora di medie dimensioni (3000 linee di testo libero)
- Dominio troppo piccolo per convergenza
- Alcuni risultati comunque interessanti: categorie lessicali

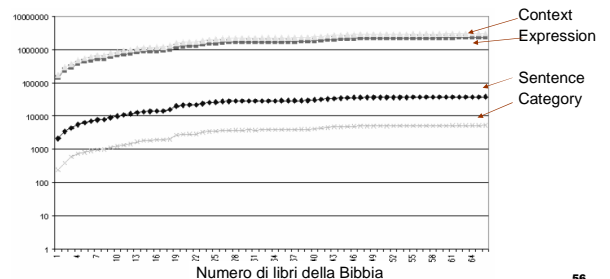
[16] ⇒ School of Medicine, University of Washington, Seattle 98195, USA	Istituti Accademici
[16] ⇒ University of Kitasato Hospital, Sagamihara, Kanagawa, Japan	
[16] ⇒ Heinrich-Heine-University, Dusseldorf, Germany	
[16] ⇒ School of Medicine, Chi a University	
[94] ⇒ Chinese	Lingue
[94] ⇒ Japanese	
[94] ⇒ Polish	
[101] ⇒ 32 : Cancer Res 1996 Oct	Riferimenti a Riviste
[101] ⇒ 35 : Genomics 1996 Aug	
[101] ⇒ 44 : Cancer Res 1995 Dec	
[101] ⇒ 50 : Cancer Res 1995 Fe	
[101] ⇒ 54 : Eur J Biochem 1994 Sep	
[101] ⇒ 58 : Cancer Res 1994 Mar	
[105] ⇒ identified in 13 cases (72	Osservazioni-verbi
[105] ⇒ detected in 9 of 87 informative cases (10	
[105] ⇒ o served in 5 (55	

16/02/2007

55

EMILE: The Bible experiment

- Testo vasto ed omogeneo (6MB)
- Il numero di espressioni e contesti incontrati, ad un certo punto, inizia a convergere



16/02/2007

56

EMILE: conclusioni

- + L'esperimento sulla Bibbia "incrina" la congettura che apprendere il linguaggio naturale usando unicamente S+ è infattibile.
- + Teoricamente, al limite, apprende la corretta struttura grammaticale del linguaggio a partire da frasi del linguaggio stesso
- + Efficiente e scalabile
 - Valido soprattutto su Large Text Corpora
- - Vengono create più categorie grammaticali del necessario
- - Non identifica i costrutti più semplici: NP VP; solo delle variazioni di essi.
- - Se parametri sui cluster ammissibili (% di completezza del blocco, nella matrice context/expressions, richiesta) sono poco ristretti, c'è rischio di sovra-generalizzazione
- Possibili applicazioni:
 - Apprendimento di informazioni semantiche da testo
 - Dominio del World Wide Web
 - Apprendimento delle strutture morfologiche delle parole

16/02/2007

57

An Evolutionary Approach (Algoritmo Genetico)

- Proposto da M. Aycinena et altri. [AYCINENA]
- Evolve una grammatica per il linguaggio naturale utilizzando un algoritmo genetico
- Supervised: il testo utilizzato è annotato con part-of-speech tags (Inglese)
- Utilizza sia S+ che S-
- Genera CFG: rappresentate come stringhe.
- Valutazione dello stato di "forma" di una Grammatica con una fitness evaluation function

16/02/2007

58

Evolutionary: Esempio di individuo

- Le Grammatiche sono Stringhe di non-terminali e pre-terminali (part of speech tags)
- Ogni 3 caratteri una produzione
SABABCBCDAE
S→AB
A→BC
B→CD
C→AE

16/02/2007

59

Evolutionary: Algoritmo

0. La popolazione iniziale (insieme delle possibili grammatiche) è generata casualmente da una distribuzione uniforme di cromosomi di una certa lunghezza
1. Esegui ciclo scegli-alleva-rimpiazza
 - A. Scegli un individuo (grammatica) a caso
 - B. **Accoppialo** con il suo compagno ideale (in base alla fitness function) in modo da generare due figli
 - Rimpiazza il genitore più debole con il figlio più in forma

16/02/2007

60

Evolutionary: Incrocio e Mutazione

- La procreazione dei due nuovi figli avviene nel modo seguente
 - Incrocio:
 - seleziona una produzione a caso in entrambi i genitori
 - Scegli un punto di partenza nelle produzioni (1° o 2°)
 - Scambia le produzioni da quel punto in poi
 - Mutazione:
 - La mutazione consiste nello scambiare un simbolo non-terminale o pre-terminale con un altro simbolo non-terminale o pre-terminale
 - Per ogni simbolo nella stringa di G, decidiamo se effettuare la mutazione

16/02/2007

61

Evolutionary: Fitness Evaluation Function

- Lo stato di forma della Grammatica è
 - in relazione diretta al Numero di frasi del Positive Set correttamente riconosciute
 - in relazione inversa al Numero di frasi generate casualmente (Negative Set) riconosciute erroneamente dalla grammatica
 - Scontato esponenzialmente, se di lunghezza inferiore

16/02/2007

62

Evolutionary: Experiments

- L' algoritmo è stato eseguito su testi letterari in lingua inglese
 - The Wizard of Oz, Alice in Wonderland, Tom Sawyer, Brown linguistic corpora
- Ogni testo è stato annotato con i seguenti sette pre-terminali, una frase per riga
 - Nouns, Verbs, adJective, adveRbs, Prepositions, deTerminers, Other

16/02/2007

63

Evolutionary: Positive and Negative Set

- Il Solo Positive Set porta ad una sovra-generalizzazione della Grammatica
 - $S \rightarrow TN | TJ | NV | NJ | PN | PV | PJ | RJ | VJ | RS | NS | PS | TS | JS | VS | SV | SJ | SN$
 - Contiene regole per appendere e pre-pendere troppi simboli. Non utilizzate nel linguaggio naturale
- Il Negative Set è dato da frasi generate casualmente: alta probabilità che non appartengano al linguaggio naturale
 - $S \rightarrow PJ | PN | TN | TJ | VJ | JJ | RJ | NJ | PS | VS | NS | JS | RS | TS | SN$
 - Consente di appendere un unico simbolo: Noun
 - In inglese, nella maggior parte dei casi, si può appendere un nome ed avere una frase sintatticamente corretta
 - "Es: "The cat" "The black cat" "The black pussy cat"

16/02/2007

64

Evolutionary: Risultati esperimenti

0→V0
0→R0
0→0J
0→VN
...

Corpus	Number of generations completed	grammar	% positive examples parsed	% negative examples parsed	fitness
aliceinwonderland	20000	0V0R00000V0N0R0J0P0P0J00 0000V0J0P00N0N000T00J000 000TJ	92.50%	8.40%	657.364
brown1_a	48500	0H44JN07N6J0J0V0P03R0N0 040N0V0R027T00T400T0V0R 0P40V0000R00J0V400N0JN	94.10%	6.10%	1967.85
brown1_b	20000	0J0N0R080N00T00N0J00V0 0V0N0P0T0T0T0P0R0J0P0J0V J	94.70%	6.70%	1227.17
brown1_c	15500	0447R47T7T1VJ22P40J0P7JN 40T7T04400140V1044P700P 708000R77JN0J7APJITR075 0N040P70N07J0P500R0N0V V0274P034N0T0M0J0N0V77J0 N44050024J0N0R02N0T0R0P 0N0T4T004P0N0P3V01154R0	80.50%	4.70%	302.998
brown1_d	4000	3V00J2V7V0V0R0J0N0N0V0T N1T33N00T3TV0N0T33N0S J00P03J0R0R0P00R030T0P0S NP	88.20%	5.00%	583.598
brown1_e	12200	0P0P00V00N00T0N0T00T0V0J 0J0R0J00R00P0N0N0N0J	93.90%	5.90%	1762.67
children	20000	0P0V0N0V0J0N0J00N0N00T0J 0V0J00R00R0J0P00T0N0T0	91.80%	5.70%	677.211
tomswayer	20000	0P0V0N0N00V0T00T0N0T0J0 N020N0P0P0P0V0J0V00R00 J00N00R0J	92.70%	8.60%	2382.89
wizardofoz	20000	0R00P0V0V00000P00R0J0V0 0T0R0000J0T00P0N0V0V0J 0TJ0N00N0N0JH	89.50%	9.20%	620.852

16/02/2007

5

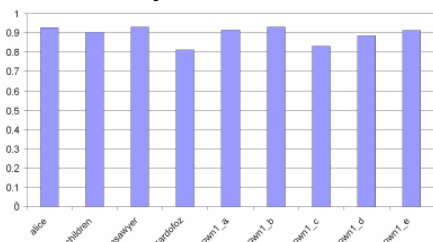
Evolutionary: risultati esperimenti

- Risultati accettabili:
 - G riconosce la maggior parte di S+
 - e una piccola parte di S-
- Col tempo, tende a ridurre la lunghezza della Grammatica
- Per valutare meglio i risultati è stata eseguita una cross-validation
 - Comportamento delle Grammatiche su testi diversi da quelli usati per la generazione.

16/02/2007

66

Evolutionary: cross-validation



Il Metodo si è dimostrato efficace

16/02/2007

67

Evolutionary: Conclusioni

- + Produce Grammatiche con alta precisione
- - La Grammatica indotta è molto diversa da quella costruita intuitivamente dall'uomo
- In Particolare, La G per l'inglese, accetta ogni stringa che termina con uno dei bigrammi previsti nelle produzioni, seguito da zero o più nomi
 - + Riconosce abilmente frasi valide
 - + Rifiuta la maggior parte delle frasi che non sono inglesi
 - - Non ci dice niente sulla struttura delle frasi
 - - Comunque, accetta troppe frasi che non sono valide

16/02/2007

68

Sommario

- Introduzione
- Grammar Induction: Teorie e Metodi
- Tecniche di apprendimento su Linguaggi Formali
- Tecniche di apprendimento sul Linguaggio Naturale
- **Risultati e Conclusioni**
- Riferimenti

16/02/2007

69

Risultati e Conclusioni

- Costruire grammatiche manualmente è un lavoro estremamente lungo ed inoltre non da garanzie di precisione sui risultati
- La Grammar Induction può essere di grande aiuto in questo processo
 - Con metodi completamente automatici (Unsupervised)
 - EMILE
 - Applicazioni a testi con impossibilità di annotazione
 - Con metodi Semi-Automatici [Chung]
 - Possibilità di dare in pasto all'algoritmo conoscenze aggiuntive
 - Possibilità di post-processare la Grammatica indotta manualmente
 - Risultati confortanti
 - Con metodi che prevedono una precedente annotazione del testo
 - Evolutionary approach
 - Risultati precisi

16/02/2007

70

Risultati e Conclusioni (2)

- Applicazione utile nei processi di annotazione del testo automatizzati
 - Metodi Supervised utilizzabili in pratica
- Alcune applicazioni utili, soprattutto per i linguisti
 - Ricerca di informazioni sulla struttura sintattica delle frasi
 - Riconoscimento di categorie semantiche
- **Nessuno pretende di modellare l'acquisizione del linguaggio umano con l'induzione di CFG**
- Grammatiche Context-Sensitive per risultati più vicini al linguaggio umano
 - ADIOS model (Automatic Distillation of structure) [Solan] – metodo unsupervised sui ampi testi non-annotati

16/02/2007

71

References

- [Gold, 1967] Gold, E. M. (1967) Language identification in the limit. *Information and Control*, 10(5):447--474.
- [Oncina, 1992] J. Oncina, P. García. Identifying regular languages in polynomial time. *Advances in Structural and Syntactic Pattern Recognition World Scientific*, 1992 99-108
- [Lankhorst, 1994] Grammatical inference with a genetic algorithm. In *Proceedings of the 1994 EUROSIM Conference on Massively Parallel Applications and Development*, pages 423-430
- [Honkela, 1995] T. Honkela, V. Pulkki, T. Kohonen. Contextual relations of words in grimm tales, analyzed by self-organizing map. In *Proceedings of the International Conference on Artificial Neural Networks*.

16/02/2007

72

References (2)

- [Baker, 1979] J. K. Baker Trainable grammars for speech recognition. In *Speech Communication Papers for the Ninety-seventh Meeting of the Acoustical Society of America*, pages 547-550
- [Lari & Young, 1990] K. Lari and S. J. Young. 1990. The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer Speech and Language*, 4:35-56
- [Wolff, 1988] J. G. Wolff. 1988. Learning syntax and meanings through optimization and distributional analysis. In *Categories and Processes in Language Acquisition*, pages 85-98
- [Adriaans, 2000] P. Adriaans, M. Trautwein, M. Vervoort. Towards High Speed Grammar Induction on Large Text Corpora

16/02/2007

73

References (3)

- [Adriaans, 1992] P. Adriaans. Language learning from a categorial perspective. Ph.D. thesis.
- [Adriaans, 1999] P. Adriaans. Learning shallow context-free languages under simple distributions.
- [Chung] Semi-Automatic acquisition of domain-specific semantic structures
- [AYCINENA] An Evolutionary approach to natural language grammar induction
- [SOLAN] Unsupervised context-sensitive language acquisition from a large corpus

16/02/2007

74