

# Testi, linguistica e ingegneria della lingua

1

LINGUISTICA COMPUTAZIONALE E  
INFORMATION TECHNOLOGY

ISABELLA CHIARI  
UNIVERSITÀ LA SAPIENZA DI ROMA

seminario DIBE, Università di Genova, 4/06/07

## *La disciplina e il suo oggetto*

2

Nascita

- fondazione dell'Association of Computational Linguistics (ACL) nel 1962

Fisionomia

- pluralità di programmi di ricerca e metodologie
- interdisciplinarietà e multidisciplinarietà

Obiettivi

- applicazioni destinate a specialisti del linguaggio
- applicazioni informatiche di uso comune

seminario DIBE, Università di Genova, 4/06/07

## Le direzioni della linguistica computazionale 1

3

### NLP (Natural Language Processing)

- Generazione del linguaggio
- Analisi del linguaggio naturale
- Correttori ortografici e grammaticali

### MT (Machine-Translation)

- Traduzione automatica
- Strumenti di ausilio (database, lessici, tesauri)

### Tecnologie del parlato

- TTS (text-to-speech)
- Speech recognition
- Dialogue systems
- Parlato multimodale

seminario DIBE, Università di Genova, 4/06/07

## Le direzioni della linguistica computazionale 2

4

### Linguistica dei corpora

- Analisi testuale/statistica linguistica
- Analisi dei dati testuali
- Integrazione in applicazioni corpus-based

### Glottodidattica

- CALL (*Computer assisted language learning*)
- *Testing* computerizzato

### Lessicografia computazionale

- Dizionari informatizzati
- Dizionari-macchina

### Semantic Web

seminario DIBE, Università di Genova, 4/06/07

## Settori di confine

5

### L'indicizzazione automatica

- serve a produrre delle analisi rapide dei testi raccolti (per esempio sul web) attraverso la individuazione delle parole-chiave di un testo (*keyword extraction*)

### L'information retrieval

- permette di individuare, dato un insieme ampio di documenti, solo quei documenti che soddisfano i nostri criteri di ricerca

### L'information extraction

- permette di cercare e selezionare i contenuti dei documenti appartenenti a un insieme

### Il text mining

- categorizzazione e classificazione dei documenti, la tematizzazione, l'estrazione di relazioni tra dati e il suo riversamento sotto forma di database

### La text summarization

- consente di generare automaticamente riassunti di testi, rapporti estratti da dati strutturati e testi che estraggano informazioni rilevanti o pertinenti a partire da una base dati testuale

seminario DIBE, Università di Genova, 4/06/07

## Se io interrogo una base testuale

6

### Voglio che tenga conto dei sinonimi

- Se cerco "linguistica", deve trovare anche "glottologia"

### Voglio che tenga conto delle flessioni

- Se cerco "ossidare", voglio che trovi anche "è stato ossidato", "ossida", ecc.

### Voglio che non confonda le espressioni semplici da quelle complesse

- Se cerco "vedere", non devo trovare "vedere rosso"

### Voglio che non confonda gli omografi

- "porta" di *lui porta* con *la porta è chiusa*

### Voglio che riconosca i nomi propri

- che distingua *Salvo è uscito* da *Non ti salvo*.

seminario DIBE, Università di Genova, 4/06/07

# L'apporto della linguistica dei corpora e testuale

7

## IL TESTO E I SUOI PROBLEMI

seminario DIBE, Università di Genova, 4/06/07

## *Corpus (plur. corpora)*

8

«raccolta completa e ordinata di scritti, di uno o più autori, riguardanti una certa materia» (De Mauro, GRADIT)

«campione di una lingua preso in esame nella descrizione di una lingua» (De Mauro, GRADIT)

TESTI

- opere di Alessandro Manzoni
- lettere d'amore,
- atti giudiziari
- perizie psichiatriche
- testi di telefonate

SCOPI

- usare le osservazioni condotte su un corpus campionario per estenderle all'intera popolazione
- comparare le osservazioni condotte su diversi corpora per confrontarle infine con un corpus di riferimento, individuandone le deviazioni

seminario DIBE, Università di Genova, 4/06/07

## I testi sono caratterizzati da...

9

### Creatività linguistica

- in creatività basata su regole (*rule-based creativity*) e creatività che cambia le regole (*rule-changing creativity*) (N. Chomsky)
- Ogni testo (“occorrenza comunicativa”) è unico e irripetibile
- Atto individuale di *parole* (F. de Saussure)

### Regolarità e prevedibilità

- una *regola linguistica* può essere vista come la descrizione di una pratica linguistica.
- la regola rappresenta una semplice regolarità, una tendenza che preferisce determinate soluzioni in una lingua, rispetto ad altre meno frequenti

seminario DIBE, Università di Genova, 4/06/07

## Pluralità

10

# Linguistica quantitativa

Approcci  
logico-  
matematici

Mirano a fornire modelli  
matematici del  
funzionamento delle lingue

Approcci  
statistici

Mirano all'estrazione di  
regolarità statistiche da  
grandi quantità di raccolte  
testuali

Approcci di  
tipo psico-  
linguistico

Intendono sottolineare il  
ruolo dei processi  
probabilistici  
nell'apprendimento, nella  
produzione e nella ricezione  
linguistica

seminario DIBE, Università di Genova, 4/06/07

## 1. Approccio di tipo logico-matematico

11

### Linguistica matematica

- ha obiettivi di tipo modellistico e predittivo
- individua modelli e rappresentazioni matematiche delle strutture linguistiche a diversi livelli
- usa strumenti di tipo algebrico

### Esponenti (a diversi livelli)

- Solomon Marcus
- Igor Mel'chuk
- Zellig S. Harris
- Noam Chomsky
- Maurice Gross

seminario DIBE, Università di Genova, 4/06/07

## 2. Approccio di tipo statistico

12

### Statistica linguistica o linguistica probabilistica

- ha obiettivi di tipo statistico-descrittivo
- presta particolare attenzione al lessico e maggiore attenzione alle realtà testuali

### Esponenti principali

- George K. Zipf
- Benoit Mandelbrot
- Pierre Guiraud
- Charles Muller
- Gustav Herdan

seminario DIBE, Università di Genova, 4/06/07

## Il lessico nei testi

13

Le parole si distribuiscono in modo regolare

- I vocabolari di base: FO, AD, AU

È possibile individuare le parole chiave di un testo

- Attraverso il confronto con corpora di riferimento
- Attraverso indici interni al testo (TFIDF)

Le parole grammaticali costituiscono il profilo del testo

Ricchezza lessicale e leggibilità

seminario DIBE, Università di Genova, 4/06/07

## 3. Approccio di tipo psicolinguistico

14

**Sottolinea come l'interiorizzazione dei fattori statistici giochi un ruolo:**

- nella performance linguistica a livello sia di produzione sia di comprensione
- nella fonologia, fonotassi
- nell'accesso al lessico
- nei meccanismi di lettura e scrittura
- nell'apprendimento della lingua materna e delle seconde lingue

seminario DIBE, Università di Genova, 4/06/07

# Le caratteristiche incalcolabili delle lingue naturali

15

PANORAMICA LINGUISTICA

seminario DIBE, Università di Genova, 4/06/07

## *La potenziale infinitezza dei segni*

16

- Le lingue storico-naturali, così come i calcoli, possono produrre un numero potenzialmente infinito di segni
- Posso creare sempre nuove frasi, e posso creare nuovi lessemi, nuove parole, che esprimano nuovi significati
- Non vi è limite di lunghezza nella produzione dei segni
- L'inventario delle unità di prima articolazione (i morfi, dotati di significante e significato) è aperto

seminario DIBE, Università di Genova, 4/06/07



### *Le famiglie di sensi e l'estensibilità dei significati*

17

- I significati sono infatti organizzati al loro interno in *accezioni*
  - ✦ polarizzazioni dei sensi in famiglie
- *Calcio*
  - “colpo dato con il piede”
  - “sport”
- In inglese, *kick, soccer, football...*
- *Estensibilità*, ossia la capacità nel tempo di sviluppare nuovi usi per rispondere ai bisogni comunicativi
- *Navigare sul web*
  - *surf the net*

seminario DIBE, Università di Genova, 4/06/07

### *Le sinonimie*

18

- Parziali sovrapposibilità di possibili sensi in alcuni enunciati prodotti o producibili
- Si istituiscono sul piano della *parole*.
- Incalcolabilità delle sinonimie
  - *dizionario*
  - *vocabolario*
    - ✦ *vai a controllare nel dizionario/vocabolario il significato di “obsolescente”*
    - ✦ *nel vocabolario di Gadda si trovano numerosi dialettalismi*

seminario DIBE, Università di Genova, 4/06/07

## WordNet <http://wordnet.princeton.edu/>

19

Elaborato al Cognitive Science Laboratory dell'Università di Princeton, ideato da G. A. Miller

Repertorio lessicale della lingua inglese organizzato per insiemi semantici

Gruppi di significati e delle gerarchie semantiche

Applicazioni di WordNet

- identificazione delle accezioni delle parole
- *information retrieval*
- identificazione delle collocazioni
- gestione di terminologie
- disambiguazione semantica
- sviluppo di ontologie

Per la lingua italiana MultiWordNet

- <http://multiwordnet.itc.it>

seminario DIBE, Università di Genova, 4/06/07

## MultiWordNet

20

Il dizionario è essenzialmente fondato sulle relazioni di sinonimia

Insiemi sinonimici, detti *synsets*

- {elaboratore, computer, cervello\_elettronico, calcolatore}

Ha censito 58.000 sensi della lingua italiana, e individuato 32.700 *synsets*

*has\_hyponym* {macchina}

*has\_hyponym* {calcolatore\_analogico}, {calcolatore\_digitale}, ecc.

*has\_part* {microchip, chip}, ecc.

{elaboratore, computer, cervello\_elettronico, calcolatore}

*corrisponde a*

{computer, data\_processor, electronic\_computer, information\_processing\_system}

seminario DIBE, Università di Genova, 4/06/07

## Le omonimie assolute e testuali

21

- Gli omonimi, infatti, sono parole caratterizzate da un significante comune, ma che rimandano a significati radicalmente diversi, spesso senza alcuna parentela etimologica
- *Omonimi assoluti*
  - *Calcio*
    - ✖ “pedata”
    - ✖ “Ca”
    - ✖ “impugnatura di un fucile o pistola”
- *Omonimi testuali*
  - *Faccia*
    - ✖ “viso”
    - ✖ “voce del verbo fare”

seminario DIBE, Università di Genova, 4/06/07

## Strumenti

22

### Omografi assoluti

- *Word sense disambiguation*
- Strumenti probabilistici/statistici

### Omografi testuali (relativi)

- POS tagging e Lemmatizzatori
  - Basati su regole
  - Probabilistici
- Operazione di base per ogni corpus

seminario DIBE, Università di Genova, 4/06/07

## Se non risolvo il problema ottengo...

23

Forma grafica	Occorrenze totali	Lunghezza	CAT	CAT_AC	CAT_SEM	Imprinting
di	12.906	02				
e	6.416	01				
in	4.221	02				
che	4.063	03				
un	3.937	02				
la	3.566	02				
a	3.218	01				
il	3.175	02				
è	2.973	01				
una	2.562	03				
per	2.557	03				
del	2.049	03				
i	1.896	01				
l	1.697	01				
da	1.649	02				
si	1.629	02				
della	1.545	05				
dei	1.533	03				

## E se guardo le concordanze trovo:

24

Intorno sinistro	Forma grafica	Intorno destro
opac . sbn . it / . Una ricerca di questo tipo ci	porta	, in un ´ ora circa di lavoro dalla scrivania di casa
RealOne Player , è in grado di collegarsi tramite la	porta	USB o seriale del PC con alcuni riproduttori portatili
parte dei casi , utilizzare quella , collegandola alla	porta	USB o firewire . Occorre però controllare il manuale
collegamento con il vostro computer alcuni utilizzano una	porta	USB , altri richiedono una scheda di rete ( offrendo
dalla necessità di economizzare spazio e risorse ( che	porta	a trascurare molti newsgroup considerati ´ minori
chiave ´ lute ´ , svolta sul catalogo principale , ci	porta	a un lungo elenco di voci , ciascuna delle quali accompagnata
termini sulla stringa   Jane Austen   condotta su Google	porta	a un elenco di circa duecentoventitemila pagine disponibili
formato HTML con un apposito bottone ´ Installa ´ che	porta	a compimento in modo automatico l ´ installazione dell
su newsgroup , mentre la linguetta ´ Directory ´ ci	porta	ai già citato indice sistematico che Google mutua da
naturale esaurirsi del ´ piccolo ´ iniziale di accessi ,	porta	al ritorno prepotente dei siti informativi di riferimento
patrimonio del museo . Ciascuna di esse a sua volta	porta	all ´ elenco per secoli e nazioni delle opere , da
ancora cd . . ) sale di una directory . Ad esempio	porta	alla directory ´ pub ´ se ci si trova in ´ pub / antivirus
iniziare la lettura dalla prima pagina ; ´ Ultima letta ´	porta	alla pagina letta l ´ ultima volta che si stava usando
stava usando il libro ; ´ Pagina letta più elevata ´	porta	alla pagina più avanzata che si è letta . È possibile
all ´ algoritmo che dal nome o indirizzo universale	porta	alla stringa di localizzazione effettiva del protocollo
soffermarsi sulle seguenti : * ´ Appearance ´ :	porta	alle schede di impostazione delle caratteristiche relative
visualizzazione di pagine multilingua ; * ´ Composer ´ :	porta	alle schede di impostazione delle preferenze per l
per l ´ uso del modulo Editor ; * ´ Advanced ´ :	porta	alle fondamentali schede di configurazione relative
che emergono . . . quando un certo numero di persone	porta	avanti delle discussioni pubbliche sufficientemente
gran parte da compiere ; l ´ avvento degli ipertesti	porta	con sé problematiche finora poco esplorate , ed è probabile
informazione al suo interno . Questo vantaggio , tuttavia ,	porta	con sé un rischio non indifferente : purtroppo , infatti
´ isolamento sociale che in molti casi l ´ handicap	porta	con sé . Attraverso Internet un disabile ha infatti
allargamento per certi versi ´ naturale ´ , ma che	porta	con sé importanti conseguenze , aprendo nuovi orizzonti
informazioni : portal , portale Sito Internet che offre una ´	porta	d ´ ingresso ´ alla rete ricca di servizi per gli utenti

seminario DIBE, Università di Genova, 4/06/07

## Quanti sono in media gli omografi in italiano?

25

Tullio De Mauro in *Capire le parole* (1999) riporta:

- Il tasso di omonimia relativa o testuale è
  - Scritti tecnici (economia e finanza): 38,6%
  - LIP 46%

**Il tasso di omografia testuale dipende dalla tipologia testuale**

- Testi con parole più brevi (come le trascrizioni del parlato, ecc.) tendono ad avere più omografi dei linguaggi tecnico-specialistici
- È una conseguenza della legge di Zipf sul numero dei significati e della saturazione (Guiraud) maggiore nelle parole brevi.

seminario DIBE, Università di Genova, 4/06/07

## Alcuni esempi

26

TESTI	<i>Codice penale</i> <b>l. giuridico</b>	<i>Internet 2004</i> <b>informatica</b>	<i>Caos Calmo (S. Veronesi) narrativa</i>	<i>LIP Corpus Roma parlato</i>
<i>Token</i>	68.728	254.365	123.781	135.716
<i>Types</i>	5.160	20.130	16.819	11.922
<b><i>Omografi testuali</i></b>	<b>3.486</b>	<b>123.097</b>	<b>60.018</b>	<b>65.630</b>
<b><i>% omografi</i></b>	<b>47%</b>	<b>49%</b>	<b>49%</b>	<b>49%</b>
<i>Parole sconosciute</i>	4%	4%	11%	6%

seminario DIBE, Università di Genova, 4/06/07

## POS tagging e lemmatizzazione

27

### Il *part-of-speech* (POS) *tagging*

- etichettatura automatica per categorie grammaticali
- Il *tagger* riceve in input una frase e restituisce in output le forme grafiche delle parole accompagnate da etichette che segnalano la categoria grammaticale di appartenenza

### ESEMPIO: la forma grafica <LA>

- potrebbe corrispondere a tre etichettature grammaticali possibili:
  - **determinante** (articolo)
  - **nome** (nota musicale)
  - **pronome** (pronomi personale)

### Tipologie

- **Tagger basati su regole** (dizionario-macchina e grammatica)
- **Tagger probabilistici** (training, parametri, applicazione statistica)

seminario DIBE, Università di Genova, 4/06/07

## Un esempio: Treetagger

28

### Autori

- Helmut Schmid, Institute for Computational Linguistics of the University of Stuttgart

### Gratuito e condiviso

- Scaricabile (Mac, Windows, Linux)
- Online (max 2 mega): <http://cental.fltr.ucl.ac.be/~pat/tagger/>

### Tagger probabilistico

- Usa decision trees
- Che determina automaticamente l'ampiezza del contesto per calcolare le probabilità di transizione (più adatto delle catene markoviane per eventi rari)
- 96,36% di precisione sul Penn-Treebank (inglese)

seminario DIBE, Università di Genova, 4/06/07

## Output di treetagger

29

```

le ART la
persone NOUN personale
si ADV si
usano VER:fin usare
la ART la
loro DET:poss loro
abilità NOUN abilità
per PRE per
essere VER:infi essere
più ADV più
umani ADJ umano
allora ADV allora
vuoi VER2:fin volere
ricordarmi VER:infi:cli ricordare
nuovamente ADV:mente nuovamente
il ART il
titolo NOUN titolo
? SENT ?
allora/cerca NOUN allora/cerca
di PRE di
dare VER:infi dare
di PRE di
spiegare VER:infi spiegare
quello PRO:demo quello
che CHE che
hai AUX:fin avere
scritto VER:ppast scrivere
e CON e
questo PRO:demo questo
a PRE a

```

seminario DIBE, Università di Genova, 4/06/07

## PROBLEMI con treetagger

30

### I parametri

- non vanno bene per tutte le tipologie testuali,
  - ad esempio il parlato
- Spesso costruire un *training corpus* ah hoc non è possibile (1.000.000 di tokens, manualmente corretti)
- Il tagger va comunque sottoposto a nuovo training se si vuole ampliare il suo lessico

### Errori sistematici

- Participi e aggettivi
- Mancato riconoscimento di nomi
- Mancata indicazione di polirematiche

seminario DIBE, Università di Genova, 4/06/07

## Training e correzione manuale

31

### Training corpus

- Se si dispone già di **un ampio corpus annotato**
- Se la tipologia è molto **uniforme**, e il vocabolario è **ridotto**
  - es. meteo, oroscopo, ricette, istruzioni per l'uso, ecc.

### Correzione manuale

- Se il corpus è **piccolo**
- oppure
- Se è assolutamente **necessaria una corretta annotazione**
  - per esempio se si vuole pubblicare il corpus di uno o più testi di un autore
- Se si può contare su un numero ampio di collaboratori

seminario DIBE, Università di Genova, 4/06/07

## Le polirematiche

32

- Le *polirematiche* sono particolari espressioni composte da più di una parola grafica, che tuttavia si comportano semanticamente e spesso morfo-sintatticamente come un solo lessema
  - *stare a cuore, forza pubblica, prigioniero politico*
- «specifico sovrappiù semantico, vale a dire la non ricostruibilità del loro significato in base alla semplice somma dei significati dei singoli componenti» (De Mauro)
- cristallizzazione morfo-sintattica
  - *voi due siete proprio due occhi di lince*
  - non \**voi due siete proprio due occhi di linci*

seminario DIBE, Università di Genova, 4/06/07



## Le collocazioni (Firth)

33

- Combinazioni di parole relativamente più libere delle polirematiche, ma accomunate da una particolare frequenza d'uso, ossia dalla preferenza per l'occorrenza congiunta dei suoi componenti.
  - *compilare un modulo*
  - *obliterare il biglietto*
  - *delitto efferato*
- Gli elementi che entrano a far parte di una collocazione sono molto più rigidi e poco analitici, quindi anche i traduttori in una lingua straniera tendono a essere imprevedibili

seminario DIBE, Università di Genova, 4/06/07

## NORMALIZZAZIONE: CHE COS'È?

34

### Pre-trattamento ortografico

- La riduzione di ambiguità dovute alle convenzioni ortografiche
- individuazione un insieme di simboli come **alfabeto** (*a, b, c, 5, 8*) e un insieme di **separatori** (*.,:;/?!*)
- ogni simbolo (il punto, la virgola, la barra, ecc.) sia univoco, ossia non venga utilizzato in modi diversi nello stesso corpus

### Pre-trattamento linguistico

- Riconoscimento di strutture cristallizzate
  - come sigle, titoli, toponimi, nomi propri (prima di ridurre eventualmente le maiuscole)
- Riconoscimento di locuzioni grammaticali e polirematiche note (da lista)

seminario DIBE, Università di Genova, 4/06/07

## Un esempio: taltac 2

35

*Trattamento automatico lessicale e testuale  
per l'analisi del contenuto di un corpus*

Sergio Bolasco

- Università La Sapienza di Roma (Economia)

Analisi lessicali e testuali

- Integrazione con risorse di riferimento (vocabolari, lessici di frequenza)
- con altri programmi di trattamento statistico (Lexico, Spad)
- e linguistico

seminario DIBE, Università di Genova, 4/06/07

## Normalizzazione: definizione alfabeto

Definizione dei caratteri (alfabeto vs separatori)

<input checked="" type="checkbox"/> [TAB]	<input checked="" type="checkbox"/> 5	<input type="checkbox"/> K	<input type="checkbox"/> a	<input type="checkbox"/> w	<input type="checkbox"/>	<input type="checkbox"/> É	<input type="checkbox"/> Í	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> 6	<input type="checkbox"/> L	<input type="checkbox"/> b	<input type="checkbox"/> x	<input type="checkbox"/> Ž	<input type="checkbox"/> X	<input type="checkbox"/> Æ	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> 1	<input checked="" type="checkbox"/> 7	<input type="checkbox"/> M	<input type="checkbox"/> c	<input type="checkbox"/> y	<input type="checkbox"/>	<input type="checkbox"/> Y	<input checked="" type="checkbox"/> »	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> *	<input checked="" type="checkbox"/> 8	<input type="checkbox"/> N	<input type="checkbox"/> d	<input type="checkbox"/> z	<input type="checkbox"/>	<input type="checkbox"/> Z	<input type="checkbox"/> ¼	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> #	<input checked="" type="checkbox"/> 9	<input type="checkbox"/> O	<input type="checkbox"/> e	<input checked="" type="checkbox"/> {	<input checked="" type="checkbox"/> '	<input type="checkbox"/> G	<input type="checkbox"/> ½	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> \$	<input checked="" type="checkbox"/> :	<input type="checkbox"/> P	<input type="checkbox"/> f	<input type="checkbox"/>	<input checked="" type="checkbox"/> ~	<input type="checkbox"/> G	<input type="checkbox"/> ¾	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> %	<input checked="" type="checkbox"/> ;	<input type="checkbox"/> Q	<input type="checkbox"/> g	<input checked="" type="checkbox"/> }	<input checked="" type="checkbox"/> "	<input checked="" type="checkbox"/> ©	<input type="checkbox"/> 1	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> &	<input checked="" type="checkbox"/> <	<input type="checkbox"/> R	<input type="checkbox"/> h	<input checked="" type="checkbox"/> ~	<input checked="" type="checkbox"/> *	<input type="checkbox"/> Z	<input type="checkbox"/> 2	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> '	<input checked="" type="checkbox"/> =	<input type="checkbox"/> S	<input type="checkbox"/> i	<input type="checkbox"/>	<input type="checkbox"/> +	<input type="checkbox"/> Á	<input type="checkbox"/> 3	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> (	<input checked="" type="checkbox"/> >	<input type="checkbox"/> T	<input type="checkbox"/> j	<input checked="" type="checkbox"/> €	<input type="checkbox"/> -	<input type="checkbox"/> Á	<input type="checkbox"/> 4	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> )	<input checked="" type="checkbox"/> ?	<input type="checkbox"/> U	<input type="checkbox"/> k	<input type="checkbox"/>	<input checked="" type="checkbox"/> -	<input type="checkbox"/> Á	<input type="checkbox"/> 5	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> *	<input checked="" type="checkbox"/> @	<input type="checkbox"/> V	<input type="checkbox"/> l	<input type="checkbox"/>	<input type="checkbox"/> ~	<input checked="" type="checkbox"/> ©	<input type="checkbox"/> Á	<input type="checkbox"/> 6	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> +	<input type="checkbox"/> A	<input type="checkbox"/> W	<input type="checkbox"/> m	<input type="checkbox"/>	<input type="checkbox"/> ~	<input type="checkbox"/> Á	<input type="checkbox"/> 7	<input type="checkbox"/> Æ	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> ,	<input type="checkbox"/> B	<input type="checkbox"/> X	<input type="checkbox"/> n	<input type="checkbox"/>	<input type="checkbox"/> §	<input checked="" type="checkbox"/> ©	<input type="checkbox"/> Á	<input type="checkbox"/> 8	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> -	<input type="checkbox"/> C	<input type="checkbox"/> Y	<input type="checkbox"/> o	<input type="checkbox"/>	<input type="checkbox"/> >	<input type="checkbox"/> ±	<input type="checkbox"/> C	<input type="checkbox"/> Y	<input type="checkbox"/> ö
<input checked="" type="checkbox"/>	<input type="checkbox"/> D	<input type="checkbox"/> Z	<input type="checkbox"/> p	<input type="checkbox"/>	<input type="checkbox"/> t	<input type="checkbox"/> œ	<input type="checkbox"/> ?	<input type="checkbox"/> È	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> /	<input type="checkbox"/> E	<input checked="" type="checkbox"/> [	<input type="checkbox"/> q	<input type="checkbox"/>	<input type="checkbox"/> #	<input type="checkbox"/>	<input type="checkbox"/> ?	<input type="checkbox"/> È	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> 0	<input type="checkbox"/> F	<input type="checkbox"/> \	<input type="checkbox"/> r	<input type="checkbox"/>	<input type="checkbox"/> %	<input type="checkbox"/>	<input type="checkbox"/> ?	<input type="checkbox"/> È	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> 1	<input type="checkbox"/> G	<input checked="" type="checkbox"/> ]	<input type="checkbox"/> s	<input type="checkbox"/>	<input type="checkbox"/> %	<input type="checkbox"/> Y	<input type="checkbox"/> µ	<input type="checkbox"/> È	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> 2	<input type="checkbox"/> H	<input type="checkbox"/> ^	<input type="checkbox"/> t	<input type="checkbox"/>	<input type="checkbox"/> %	<input type="checkbox"/>	<input type="checkbox"/> ¶	<input type="checkbox"/> È	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> 3	<input type="checkbox"/> I	<input checked="" type="checkbox"/> _	<input type="checkbox"/> u	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> ¶	<input type="checkbox"/> È	<input type="checkbox"/> ö
<input checked="" type="checkbox"/> 4	<input type="checkbox"/> J	<input checked="" type="checkbox"/> `	<input type="checkbox"/> v	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> È	<input type="checkbox"/> ö
				<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		

37

## NORMALIZZAZIONE

**PUNTEGGIATURA**

- Apostrofi in accenti
- Maiuscolo/minuscolo

**POLIREMATICHE E COLLOCAZIONI (base)**

- Locuzioni gramm.
- Polirematiche nominali

**NOMI**

- nomi propri
- toponimi
- celebrità
- titoli
- Sigle

**LISTE PERSONALIZZATE**

38

## Così ottengo ad esempio...locuzioni come..

Forma grafica	Occorrenze totali	Lunghezza	CAT
quasi tutti	9	11	PRON
chissà che	1	10	PRON
tutto quanto	3	12	PRON
in tutto	11	08	PRON
tutto ciò	8	09	PRON
tutto questo	6	12	PRON
non tutti	7	09	PRON

Locuzioni preposizionali, pronominali, polirematiche, ecc.

Forma grafica	Occorrenze totali	Lunghezza	CAT
guerra fredda	2	13	N
big bang	1	08	N
punto vendita	1	13	N
Ottocento	1	09	N
corpo sociale	2	13	N
caricolo vizioso	1	15	N
agosto	1	06	N
settembre	54	09	N
personal computer	16	17	N
ministero del Tesoro	1	20	N
anni Sessanta	3	13	N
anni Settanta	1	13	N

## NOMI PROPRI, SIGLE, FORMULE

39

Forma grafica	Occorrenze totali	Lunghezza	CAT
Salvo	105	NM	
Arianna	507	NM	
Jerry	105	NM	
David	305	NM	
Alice	505	NM	
Lewis	505	NM	
Riccardo	308	NM	
Gabriel	207	NM	
George	306	NM	
Gabriele	108	NM	
Dante	605	NM	
Edward	206	NM	
Massimo	307	NM	
John	504	NM	
Alex	404	NM	

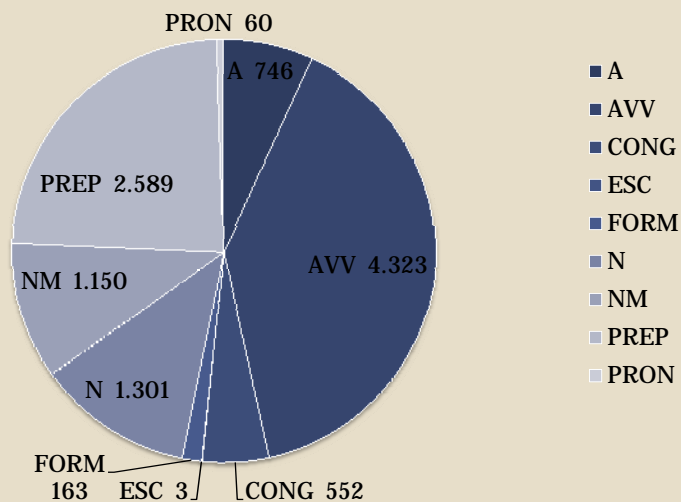
  

Forma grafica	Occorrenze totali	Lunghezza	CAT
fatto sta che	113	FORM	
non è un caso che	417	FORM	
il problema è che	117	FORM	
è un problema che	117	FORM	
non è detto che	215	FORM	
se è vero che	713	FORM	
l'idea che	110	FORM	
non è che	109	FORM	
è il caso	509	FORM	
è vero che	210	FORM	
è bene che	110	FORM	
si dice che	211	FORM	
è lo stesso	411	FORM	
è chiaro che	312	FORM	
dire che	3608	FORM	
pensare che	811	FORM	
è certo	607	FORM	
penso che	109	FORM	
sembra che	210	FORM	
sapere che	810	FORM	
è naturale	110	FORM	
per così dire	413	FORM	

seminario DIBE, Università di Genova, 4/06/07

## normalizzazione: testo internet 2004

40



seminario DIBE, Università di Genova, 4/06/07

## Prima e dopo la normalizzazione

41

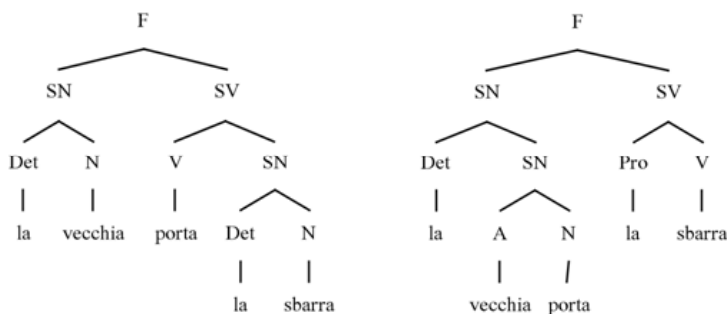
	Prima della normalizzazione		
Dati corpus	Normalizzato	Differenza	
<b>TOKEN (occorrenze)</b>	<b>254.365</b>	<b>240.173</b>	<b>14.192</b>
<b>TYPES</b>	<b>20.130</b>	<b>18.730</b>	<b>1.400</b>
<b>OMOGRAFI</b>	<b>123.097</b>	<b>108.760</b>	<b>14.337</b>
	<b>(48,4%)</b>	<b>(45,3%)</b>	

seminario DIBE, Università di Genova, 4/06/07

## Le ambiguità sintattiche

42

- ▶ Alcune frasi – prese in isolamento – possono avere diverse interpretazioni sintattiche che possono essere associate a più di un plausibile albero sintattico
  - ▶ *la vecchia porta la sbarra*
  - ▶ *lo studente ha risolto i suoi problemi col computer*



seminario DIBE, Università di Genova, 4/06/07

## *Deittici e quantificatori*

43

- **Uso dei deittici**
- **Solitamente è il contesto a permetterci di risolvere il dubbio**
  - ✦ *ho preso il libro e l'ho aperto*, in cui il pronome anaforico *lo* si riferisce a *libro*
- **Alcuni usi dei *quantificatori***
  - *tutti gli studenti danno un esame*

seminario DIBE, Università di Genova, 4/06/07

## *La metalinguisticità riflessiva*

44

- **Permette la formazione di testi nei quali ci si riferisce a elementi linguistici: la lingua è usata per parlare della lingua**
  - il pallone in inglese si dice ball
  - ho detto stilare non sfilare
- **Permette di sovvertire le normali considerazioni relative alla grammaticalità**
  - \* *con mangia le mani*
  - *con è una preposizione*

seminario DIBE, Università di Genova, 4/06/07

## *I corpora in linguistica computazionale*

45

### Lessicografia elettronica *corpus-based*

- Dizionari informatizzati
- Dizionari macchina *corpus based*

### Training corpora per il NLP

- *Taggers e parsers con training corpora*

### Traduzione automatica

- *Corpus-based*
- *Example-based machine translation*

### Tecnologie del parlato

- Addestramento allo *speech recognition*
- Sintesi *corpus-based*

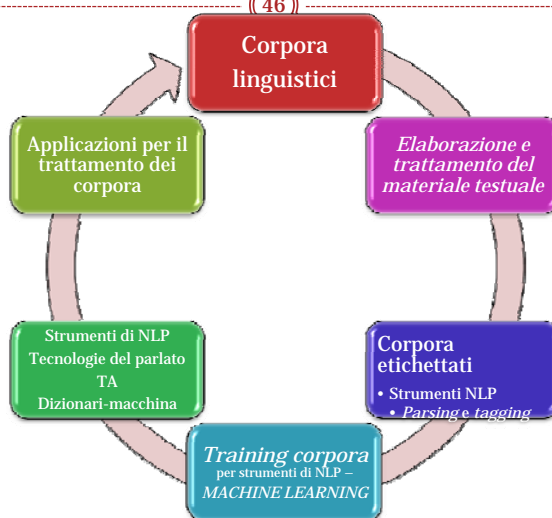
### Machine learning – Information technology

- Individuazione automatica di *patterns* estratti dai dati

seminario DIBE, Università di Genova, 4/06/07

## *Il circolo virtuoso*

46



seminario DIBE, Università di Genova, 4/06/07

**Grazie!**

47

**Le slides powerpoint**

- da martedì
- sul sito:

**[www.alphabit.net](http://www.alphabit.net)**

- sotto la voce CONVEGNI e NOVITA'

**Isabella Chiari**

- **[Isabella.chiari@uniroma1.it](mailto:Isabella.chiari@uniroma1.it)**