

# Metodi linguistici di analisi di testi

Leggere i corpora  
(2023/2024)

Giovanni Urraci

[giovanni.urraci@unipd.it](mailto:giovanni.urraci@unipd.it)



# ALCUNI STRUMENTI



**AntConc** [<https://www.laurenceanthony.net/software/antconc>]



**Voyant tools** [<https://voyant-tools.org>]



**NoSketch Engine e Sketch Engine** [<https://www.sketchengine.eu>]

**Tree Tagger** [<https://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger>]



**READ-IT** [[https://www.ilc.cnr.it/dylanlab/apps/texttools/?tt\\_user=guest](https://www.ilc.cnr.it/dylanlab/apps/texttools/?tt_user=guest)]

# ESPLORAZIONE DEI CORPORA

Linguistica computazionale, analisi dei dati testuali,  
linguistica dei corpora, linguistica quantitativa...



Etichette diverse, descrivono la stessa tipologia di approccio e  
muovono dallo stesso oggetto di studio  il **corpus**.

Consentono di maneggiare efficacemente corpora di grandi dimensioni,  
e di porre ai testi domande originali.

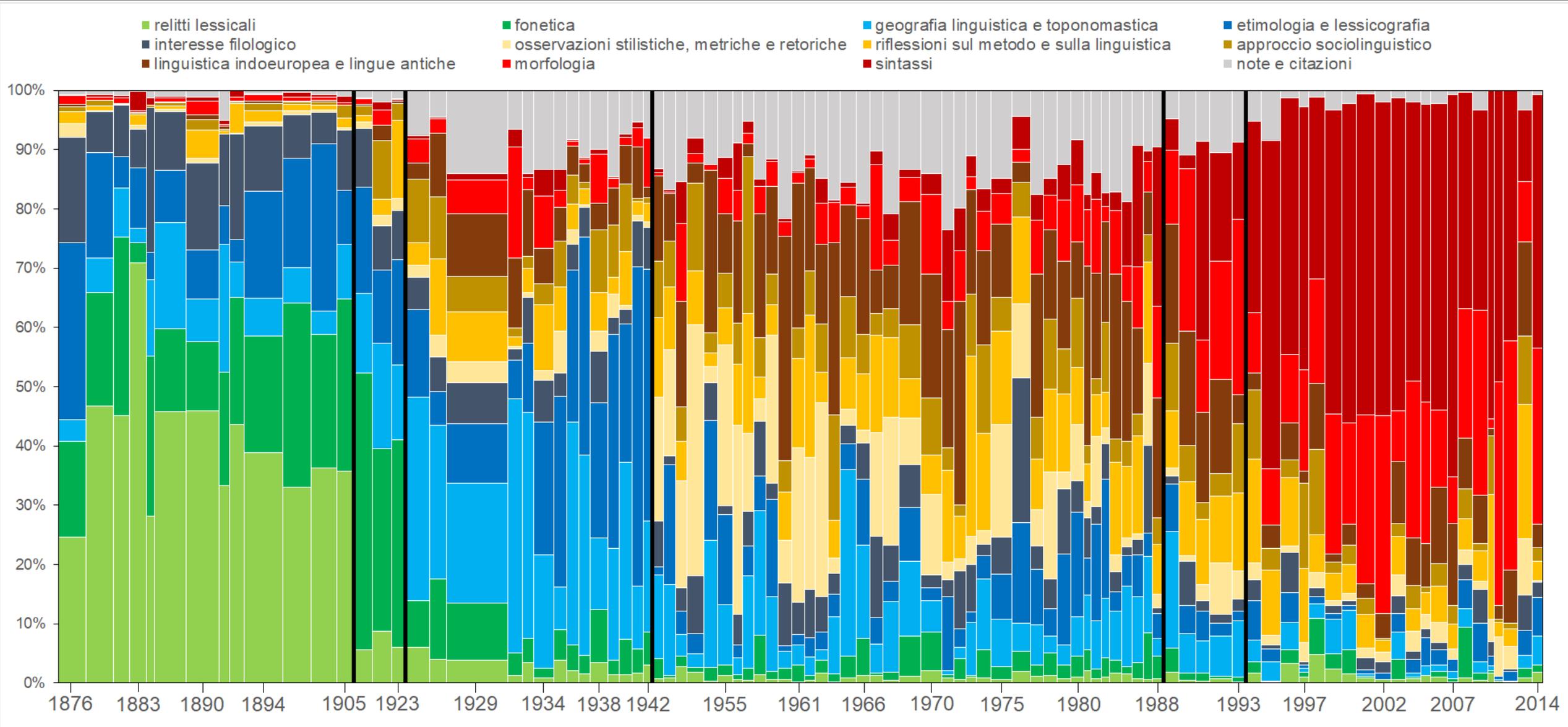
Tuttavia, non costituiscono un'alternativa ai tradizionali approcci qualitativi



**complementarietà**

# ANALISI STATISTICA DEI DATI TESTUALI

Profilo tematico dell'«Archivio Glottologico Italiano». Topic extraction, metodo di Reinert



**Corpus** raccolta omogenea di testi selezionati per la loro rappresentatività (di una lingua, una varietà, un dominio, una tipologia testuale). Oggi in formato digitale.

Possono essere di lingua scritta o trascritta, sincronici o diacronici, specialistici o 'di riferimento', di apprendenti, eccetera.

Es.: discorsi di insediamento dei presidenti della repubblica italiana; articoli pubblicati su la Repubblica tra il 1985 e il 2000.

**Sub corpus** porzione di un corpus: singolo testo o insieme dei testi che condividono una o più proprietà pertinenti ai criteri oggetto di analisi.

Indispensabili → l'analisi quantitativa si fonda sulla comparazione, necessaria per identificare le caratteristiche di uno specifico gruppo.

Es.: discorso di insediamento del presidente Sergio Mattarella; articoli di materia economica pubblicati su la Repubblica tra il 1985 e il 2000.

# ALCUNI CORPORA

- KIParla – corpus di parlato spontaneo
  - UniverS-ITA – testi formali scritti da studenti universitari
  - la Repubblica (1985-2000)
  - OpenSubtitles (2018) – raccolta di sottotitoli
  - itTenTen20
  - itWaC
- Siti Internet e testi raccolti online

# CREAZIONE CORPUS

La costruzione di un corpus varia a seconda delle variabili coinvolte e del software utilizzato.

Soluzione più frequente: testi in formato .txt (UTF-8), un documento per ogni sub corpus, raccolti in una cartella.

\*\*\*\*001 \*Rivista=AGI \*Anno=1876

Postille etimologiche.

Saggio di un Glossario Modenese ossia studii del conte Giovanni Galvani intorno le probabili origini di alquanti idiotismi della città di Modena e del suo contado. Scrisi le seguenti postille etimologiche quattro e più anni sono; e le scrissi principalmente coll'intento di mettere per così dire a fronte due scuole, la vecchia e la nuova, la scuola senza metodo e quella del metodo. Attendendo per debito d'ufficio ad insegnar glottologia nell'Ateneo torinese, mi parve che dalla pubblicazione del Galvani venissemi non solo buona occasione, ma obbligo di dimostrare come nelle cose della linguistica più non valgano gran fatto di per sé soli né ingegno, né dottrina, né squisita coltura di lettere; pregi che ni uno avrebbe potuto negare al Galvani; ma si debba innanzi tutto chiedere a quella, che ora può dirsi ed è veramente scienza delle lingue, il metodo e i principj.

[...]

\*\*\*\*002 \*Rivista=AGI \*Anno=1877

Fonetica del dialetto di Val Soana (canavese).

Il dialetto di Val Soana è parlato dalla popolazione dei quattro comuni della valle di questo nome, che sono Ingrìa, Ronco, Valprato e Campiglia. È inoltre parlato nei due comuni di Ribordone e Frassinetto, il primo de' quali sta a destra, l'altro a sinistra della valle. La popolazione di fatto ascendeva per cotesti comuni, il 31 dicembre del 1871, a 7582 anime, distribuite come segue: . Siccome però il censimento si faceva appunto in quel tempo dell'anno in cui la popolazione virile suole emigrare dalla valle, conviene aggiungere a questa cifra, per ottenere approssimativamente lo stato della popolazione di diritto, o meglio della vera popolazione effettiva, poco meno d'un altro migliajo d'anime; così che il numero delle persone, che ha per favella materna il dialetto valsoanino, riesce all'incirca di 8,500.

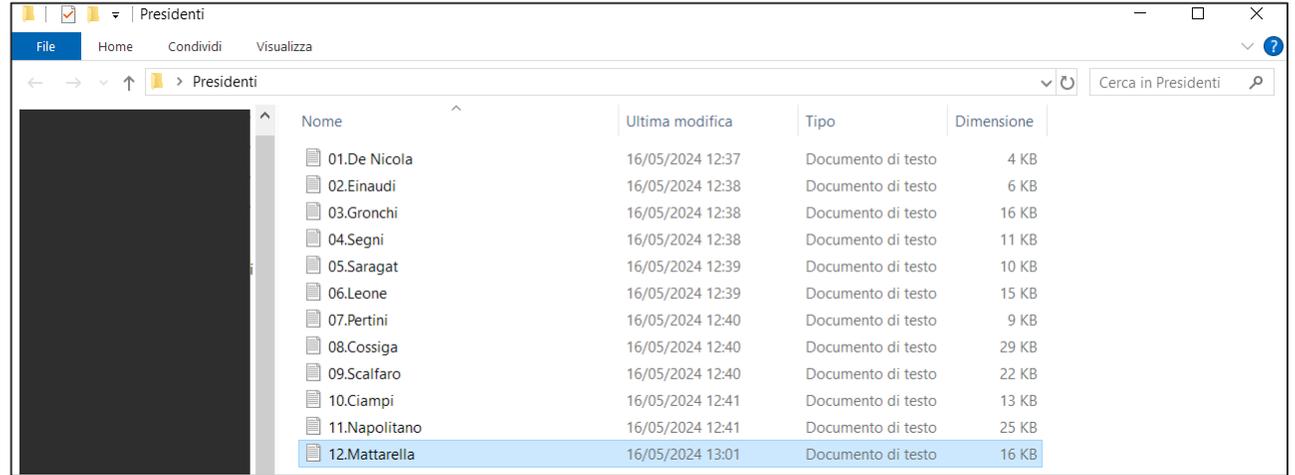
[...]

\*\*\*\*096 \*Rivista=LN \*Anno=1939

Correnti dotte e correnti popolari nella lingua italiana.

In quell'operetta che si può considerare il più antico dei lessici italiani, il Vocabulario di cinque mila Vocabuli Toschi non men oscuri che utili e necessari del Furioso, Boccaccio, Petrarca e Dante, pubblicato a Napoli nel 1536, l'autore, Fabricio Luna, cita l'aneddoto di un gentiluomo che, volendo ordinare ai suoi staffieri di accorciarli la coreggia delle staffe, si esprimeva così: O famuli, o famuli, abbreviatimi questi sustentacoli, che son troppo prolissi! E nello stesso anno un bolognese, Giovanni Filoteo Achillino, nelle Annotazioni della volgar lingua, citava così le parole del gentiluomo: Agricola, abbreviammi esto sustentaculo ch'è nimio prolisso. Dalle origini fino ad oggi, nella pratica e nella teoria, in modi diversi eppure con grandissime conformità, affiorano le proteste contro la lingua troppo dotta; e ci permettono di cogliere al vivo, grazie appunto alle esagerazioni della polemica, uno degli aspetti perpetui dell'italiano (che del resto riappare, mutati i termini, per tutte le lingue colte d'Europa), il confluire nella lingua di correnti popolari e di correnti dotte.

[...]



	A1			f2	ID	
	A	B	C	D	E	
1	ID	Day	Month	Year	Representative	Tweet
2	1	31	12	2013	Renzi	In Palazzo Vecchio, al lavoro per preparare la Giunta di fine anno. Si annuncia bella corpora...
3	2	31	12	2013	Renzi	@Fiorello Nessuno è perfetto, Fiorel! Auguri
4	3	31	12	2013	Renzi	Ho mantenuto l'impegno di andare nella terra dei fuochi. In silenzio e senza dichiarazioni. Graz...
5	4	30	12	2013	Renzi	Oggi a #Firenze, Eataly ha aperto nella via Martelli pedonale. Sono 122 posti di lavoro. E io son...
6	5	28	12	2013	Renzi	Grazie a tutti, buona serata. Torneremo con #matteorisponde dopo Natale
7	6	27	12	2013	Renzi	Pronto per il #matteorisponde Tra cinque minuti si parte...
8	7	27	12	2013	Renzi	Colonna sonora di domani al @pdnetwork: "Resta ribelle" dei Negrita "la tua canzone". Piace?
9	8	27	12	2013	Renzi	Insieme a @bobogiac nel giorno in cui finisce lo sciopero della fame. pic.twitter.com/tsNg6hOz
10	9	20	12	2013	Renzi	Caro @beppe_grillo ti rispondo nei prossimi giorni con una #sorpresina che ti sto preparando...
11	10	17	12	2013	Renzi	Scusate il ritardo nelle risposte, sono stati giorni intensi. E bellissimi. Ma il meglio deve ancora...
12	11	17	12	2013	Renzi	Grazie.
13	12	17	12	2013	Renzi	Giornata difficile da dimenticare... Ci vediamo alle 22 all'ObiHall (Firenze Sud) e in streaming s...
14	13	14	12	2013	Renzi	Grazie a tutti i volontari che consentono le primarie e ai cittadini che stanno votando. Buon vot...
15	14	14	12	2013	Renzi	La piazza di Empoli mi resterà nel cuore a lungo. Grazie ragazzi, è davvero #lavollabuona
16	15	11	12	2013	Renzi	Grazie a tutte le volontarie e i volontari che stanno prendendo freddo ai tavolini nelle mille piazz...
17	16	11	12	2013	Renzi	Mamma mia, quanto entusiasmo. Grazie a tutti. Adesso inizia la parte difficile: uno per uno, ca...
18	17	11	12	2013	Renzi	Ultimo miglio e poi si #cambiaverso Arriviamo a Milano. Questa è la #vollabuona, ora o mai più
19	18	8	12	2013	Renzi	Tra qualche minuto a Tortona con Realacci, Bastioli, Angelantoni, Ghisolfi con le proposte su a...
20	19	8	12	2013	Renzi	Un gigante. #Mandela. Ciao #Mediba



# CONCORDANZE

Voyant Tools			
Contesti			
Documento	Sinistra	Parola	Destra
01.De Nicola	la mia breve ma intensa missione di Capo provvisorio dello	stato	inspirandomi ad un solo ideale: di servire con fedeltà e
01.De Nicola	nel quale i diritti dei cittadini e i poteri dello	stato	siano egualmente garantiti, trarrà dal passato salutarì insegnamenti, consacrerà per
02.Einaudi	trapasso avvenuto il 2 giugno dall'una all'altra forma istituzionale dello	stato	fu non solo meraviglioso per la maniera legale, pacifica del
02.Einaudi	ceto politico migliore di quello che, all'alba del Risorgimento, era	stato	fornito dal suffragio ristretto. Or qui si palesa il grande
02.Einaudi	è garanzia della libertà della persona umana contro l'onnipotenza dello	stato	e la prepotenza privata; e garantire a tutti, qualunque siano
03.Gronchi	indicazione del Parlamento il mandato per la suprema magistratura dello	stato	. Quest'attesa e quest'ansia non mi sgomentano, ma fanno più presente
03.Gronchi	quali che siano le qualità benevolmente attribuitemi, l'elemento determinante dello	stato	d'animo comune a tanta parte del popolo italiano. Ma è
03.Gronchi	I problemi fondamentali del passato sono stati la ricostituzione dello	stato	nella sua organizzazione e nella sua autorità, la ricostruzione economica
03.Gronchi	sistema produttivo; ma essi hanno già nella organizzazione politica dello	stato	moderno un'influenza che è adeguata alla loro importanza economica. Io
03.Gronchi	il suffragio universale ha condotto sino alle soglie dell'edificio dello	stato	senza introdurle effettivamente dove si esercita la direzione politica di
03.Gronchi	insieme opera di progresso e di conservazione, di intervento dello	stato	e di rispetto dell'iniziativa privata. Nuove forme di organizzazione economica
03.Gronchi	e sud che travagliano penosamente l'efficienza operativa dell'economia nazionale. Lo	stato	può dare un valido concorso a nuove forme di rapporti
03.Gronchi	della nostra Costituzione. Distinguere in questo campo la responsabilità dello	stato	da quella dell'iniziativa privata non vuol dire contrapporre le due
03.Gronchi	vigile ardimento, la trasformazioni delle strutture economiche e sociali. Allo	stato	spetta in primissima istanza la responsabilità di mantenere le condizioni
03.Gronchi	le condizioni necessarie all'ordinato sviluppo democratico della comunità nazionale. Lo	stato	è imparziale tutore dei diritti di ciascuno, della libertà, dell'uguaglianza
03.Gronchi	di un'ordinata convivenza. Non è una definizione puramente giuridica lo	stato	di diritto: è l'espressione di un'esigenza politica e sociale alla
03.Gronchi	sospetta che, dietro la facciata di quel che si chiama	stato	, si cela il giuoco di potenti gruppi organizzati". Nessuna parola
03.Gronchi	definire il carattere e le responsabilità di un Capo di	stato	per quanto riguarda la piena osservanza della Costituzione, delle norme
03.Gronchi	sinceramente accettato di soggezione alla legge ed all'imparziale autorità dello	stato	. Onorevoli deputati, onorevoli senatori! La nuova fase della nostra vita
04.Segni	piena e di imparzialità assoluta nell'esercizio della suprema magistratura dello	stato	. A Luigi Einaudi, che gli successe, va il pensiero commosso
04.Segni	a me spetta determinare gli indirizzi politici nella vita dello	stato	, prerogativa questa del governo della Repubblica e massimamente di questo
04.Segni	appropriate al bene comune. Ma a me, quale Capo dello	stato	, incombe, nell'esercizio delle mie funzioni, il dovere di tutelare l'osservanza
04.Segni	affinché sia garantita, nella forma e nello spirito dell'attività dello	stato	, l'unità civile e morale della nazione italiana, una e indivisibile
04.Segni	continuità ed unità di questa nostra Repubblica che è uno	stato	di diritto, dotato di leggi giuste ed uguali per tutti
04.Segni	equanime giustizia. Considerando quelli che sono gli organi supremi dello	stato	ed i miei doveri verso di essi, mi sia consentito
05.Saragat	dei legislatori, degli uomini di governo e dei Capi di	stato	. Ma la pace si persegue creando con tenacia e con
05.Saragat	Costituzione a uno statuto di sovranità e indipendenza accanto allo	stato	sovrano e indipendente nella sfera propria. La Repubblica democratica difende
05.Saragat	il nostro secondo Risorgimento. Le relazioni tra il Capo dello	stato	ed il governo sono fissate dalla Costituzione; e sarà nella
05.Saragat	ordinata convivenza goverà all'intero Paese. Dell'azione governativa i funzionari dello	stato	costituiscono lo strumento fondamentale. Comprendo tutto il travaglio cui questi
05.Saragat	la soluzione di essi consentirà di dare alla macchina dello	stato	quell'efficienza che i compiti odierni richiedono in particolare misura. Agli
05.Saragat	voi mi avete affidato con la più alta magistratura dello	stato	la custodia dei supremi valori della patria. Con animo commosso
06.Leone	potrà garantire la maggiore partecipazione del cittadino alla vita, dello	stato	, caratteristica essenziale della democrazia. Coerente con la linea politica di
06.Leone	nella società tecnologica contemporanea. Per quanto riguarda i rapporti tra	stato	e Chiesa, è nella Costituzione la direttrice di operare perché
06.Leone	fanno il giusto posto alla indipendenza ed alla sovranità dello	stato	e della Chiesa cattolica, ciascuno nel proprio ordine; si tratta
06.Leone	Alle forze armate, garanzia dell'indipendenza nazionale e della sovranità dello	stato	, nelle quali i nostri giovani trovano una grande scuola di
06.Leone	senso del dovere e nel sacrificio dei molti servitori dello	stato	un punto di fiducia. Va rinnovato in questa sede l'invito
06.Leone	sua alla coscienza giuridica e in una vivace custodia dello	stato	accompagnò l'opera di De Gasperi per la rinascita del Paese

# CONCORDANZE (2)

## Keyword In Context

### CONCORDANCE

Italian parliamentary debates (ParlaMint 2.1)

simple avversari • 258

8.4 per million tokens • 0.00084%

Get more space



Details

Left context KWIC Right context

1	<input type="checkbox"/>	<input type="checkbox"/>	CorsiniPaolo • ...	udico questo Governo? Un Governo di tregua; un Governo frutto di un armistizio. Del resto, la pace si stipula tra <b>avversari</b> . Tra loro lo scontro non è così cruento come tra consanguinei. Allora, come tradurre una costrizione, uno stato c
2	<input type="checkbox"/>	<input type="checkbox"/>	SustaGianluca • ...	di costruire una democrazia dell'alternanza in cui reciprocamente ci si rispetti, ci si riconosca, anche come <b>avversari</b> . Questo non passa attraverso l'annullamento delle distinzioni o la rinuncia, anche in questo periodo di conviven:
3	<input type="checkbox"/>	<input type="checkbox"/>	SchifaniRenato ...	oma , possa avvalersi di una captatio benevolentiae, di una acquisizione di consenso maggiore rispetto agli altri <b>avversari</b> in quanto parlamentare in carica e vuole dunque sfrondare il campo da questo sospetto per presentarsi ai propri
4	<input type="checkbox"/>	<input type="checkbox"/>	ScilipotilsgroD...	del premio di maggioranza alla Camera , poiché il partito o la coalizione che ha conseguito un voto in più degli <b>avversari</b> ottiene il 55 per cento dei seggi: tramite il premio di maggioranza, la coalizione guidata da Bersani ha conseguit
5	<input type="checkbox"/>	<input type="checkbox"/>	CentinaioGianMa...	na visione laica della vita. Pur non condividendo alcune sue esternazioni, spesso pungenti nei confronti dei suoi <b>avversari</b> , rispettiamo la forza d'animo e la convinzione con cui ha sempre cercato di far comprendere il suo pensiero. La
6	<input type="checkbox"/>	<input type="checkbox"/>	BondiSandro • 2...	. I partiti, attraverso questa forma di lotta politica, sono quasi sempre inflessibili nei confronti degli <b>avversari</b> politici e indulgenti verso la propria parte politica. Vero, senatore De Cristofaro ? Quando l'onorevole Vendola ,
7	<input type="checkbox"/>	<input type="checkbox"/>	BondiSandro • 2...	i limiti del rispetto per le persone, specialmente quelle come il ministro Alfano , conosciuto da tutti in quest'Aula, <b>avversari</b> e amici, per la sua integrità morale e il suo senso di rispetto delle istituzioni. Manteniamo dunque il confronto ent
8	<input type="checkbox"/>	<input type="checkbox"/>	CasiniPierFerdi...	sul rispetto che sempre si deve nutrire per i membri del Parlamento , in particolare - per chi tali li ritiene - per gli <b>avversari</b> politici, e che possa responsabilmente compromettere l'esito di questo Governo. Lo sforzo che gli italiani stanno
9	<input type="checkbox"/>	<input type="checkbox"/>	BondiSandro • 2...	capacità, a cui lei ha fatto riferimento nel suo discorso, di comprendere le ragioni degli altri, le ragioni degli <b>avversari</b> politici? Ricordo quell'epoca, onorevole Letta , l'epoca del compromesso storico (lo ricorda come me), della
10	<input type="checkbox"/>	<input type="checkbox"/>	CastaldiGianluc...	. Signor Presidente , colleghi, sento il dovere, in coerenza con l'educazione che uso nei confronti di tutti, amici e <b>avversari</b> , nonché per i rapporti cordiali che ho intrapreso con molti di voi, soprattutto i colleghi della 10a Commissione , d
11	<input type="checkbox"/>	<input type="checkbox"/>	ScilipotilsgroD...	che allora avevano assunto una posizione completamente diversa abbiano capito oggi che il rispetto degli <b>avversari</b> è qualcosa di importante e che la democrazia si fonda essenzialmente sul dialogo costruttivo con l'avversario.
12	<input type="checkbox"/>	<input type="checkbox"/>	RomaniMaurizio ...	torale, eternamente rimandata per interesse momentaneo. L'interesse è farla nel momento giusto per fregare gli <b>avversari</b> e mantenere la poltrona, non farla in modo giusto per garantire al popolo un'equa rappresentanza. Come possa
13	<input type="checkbox"/>	<input type="checkbox"/>	QuagliarielloGa...	si affronta questa materia bisogna innanzitutto leggere e meditare sulle parole e sulle argomentazioni dei propri <b>avversari</b> ed evidentemente trarne insegnamento e, laddove possibile, conforto. Cito quanto ha scritto qualche mese fa un
14	<input type="checkbox"/>	<input type="checkbox"/>	ZizzaVittorio • ...	che consiste non nella presenza di Berlusconi sulla scena istituzionale e politica (come invece tanti suoi <b>avversari</b> affermano), ma nel fatto che in vent'anni non si è riusciti ad abbattere Berlusconi per via elettorale e si tenta di
15	<input type="checkbox"/>	<input type="checkbox"/>	MinzoliniAugust...	? Innanzitutto il ricordo indelebile di un omicidio politico, di un ritorno a quel passato drammatico in cui gli <b>avversari</b> venivano liquidati con tutti i mezzi possibili, meno che con quelli leciti della competizione politica. Ma resterà
16	<input type="checkbox"/>	<input type="checkbox"/>	GasparriMaurizi...	anche oggi fare quello che fa sempre: ignorare il mandato improvvidamente ricevuto di senatore a vita. A tutti gli <b>avversari</b> politici credo che dobbiamo dire con serenità che Silvio Berlusconi continuerà a svolgere il suo ruolo politico. No
17	<input type="checkbox"/>	<input type="checkbox"/>	ToniniGiorgio • ...	approvare il provvedimento e si trovano nell'impossibilità di farlo (mi riferisco, in particolare, ai colleghi e <b>avversari</b> di Forza Italia ). Allo stesso modo, sento il bisogno di rivolgere delle scuse ai colleghi delle opposizioni che sono
18	<input type="checkbox"/>	<input type="checkbox"/>	LettaEnrico • 2...	ome il passaggio da una situazione di contrapposizione tossica tra nemici a un sistema di competizione sana tra <b>avversari</b> ; un passaggio obbligato dall'esito del voto di febbraio, ma, soprattutto, dalla necessità - che io giudico, oggi più
19	<input type="checkbox"/>	<input type="checkbox"/>	SacconiMaurizio...	per essere rimessa al piccolo gioco tattico o, ancor peggio, per essere utilizzata come ennesima clava contro gli <b>avversari</b> in una dialettica politica già viziata da molte diffidenze. Il metodo ragionevole è quello che muove dalle forze
20	<input type="checkbox"/>	<input type="checkbox"/>	CandianiStefano...	qualcosa. Penso all'infelice battuta che l'altro giorno ha fatto la senatrice Pezzopane quando, rivolta ai suoi <b>avversari</b> politici, ha detto: "sterminiamoli"; poi, si è corretta dicendo che intendeva dire: "asfaltiamoli". Magari fa

Rows per page: 20 1-20 of 258 1 / 13



## PRE-PROCESSING

- Tokenizzazione
- Normalizzazione

## CLASSIFICAZIONE DELLE ENTITÀ D'ANALISI

- Named-entity recognition
- Segmentazione in frasi e assegnazione POS
- Lemmatizzazione e stemming

## MISURE LESSICOMETRICHE ESSENZIALI

- Frequenza, rango e hapax
- Media e mediana
- TTR

## INDIVIDUAZIONE DEL LESSICO RILEVANTE

- Frequenza
- Keyword
- collocazioni

# PRE-PROCESSING

Fase di preparazione del corpus.

Determina validità delle operazioni successive e qualità dei risultati.

Impone decisioni legate allo scopo dell'indagine.

**1. Tokenizzazione** Segmentazione in token, le più piccole unità di analisi: sequenze di caratteri isolate da separatori.



La macchina deve essere istruita,  
impostazioni né banali né scontate.

Es.: in un corpus di tweet # e @ sono rilevanti.  
Oppure: cosa fare con gli apostrofi? E - ?

Distinzione fondamentale:

- **word-type**: parole diverse attestate in un corpus = graficamente distinte
- **word-token**: occorrenze

Sopra	la	panca	la	capra	canta	sotto	la	panca	la	capra	crepa	TOKEN	→ 12
1	2	3	4	5	6	7	8	9	10	11	12	TYPE	→ 7
1	2	3	2	4	5	6	2	3	2	4	7		

Fase di preparazione del corpus.

Determina validità delle operazioni successive e qualità dei risultati.

Impone decisioni legate allo scopo dell'indagine.

## 2. Normalizzazione

Uniformazione grafica e pulitura del testo. Passaggio fondamentale.

Es.: maiuscolo > minuscolo (presidente ≠ Presidente)

apostrofi e accenti (perchè > perché)

uniformazione dei composti

eventuali interventi su grafie estemporanee.



Quali variazioni sono rilevanti e quali un ostacolo all'analisi?

Cosa fare con <h> etimologica e oscillazioni ortografiche?

# CLASSIFICAZIONE DELLE UNITÀ D'ANALISI

**1. Named-entity recognition** Riconoscimento e categorizzazione semantica di unità lessicali mono- e polirematiche: persone (Sergio\_Mattarella), luoghi (Monaco\_di\_Baviera), eventi (Seconda\_Guerra\_Mondiale), organizzazioni (ONU), professioni, eccetera. Inoltre: specifiche categorie di dominio (es.: farmaci, molecole).



Eseguito attraverso il confronto con liste di riferimento.

**2. Segmentazione in frasi** shallow parsing: analisi dei costituenti e individuazione dei sintagmi  
full parsing: albero sintattico completo



assegnazione POS

**3. Lemmatizzazione** Raggruppamento delle forme flesse sotto i lemmi di riferimento.  
Operazione complessa per testi non canonici (es. reiterazione vocalica e refusi nei testi da Internet): necessario l'addestramento su corpora pertinenti annotati.  
Possibile anche lo stemming: riduzione alla radice – problematico per it.

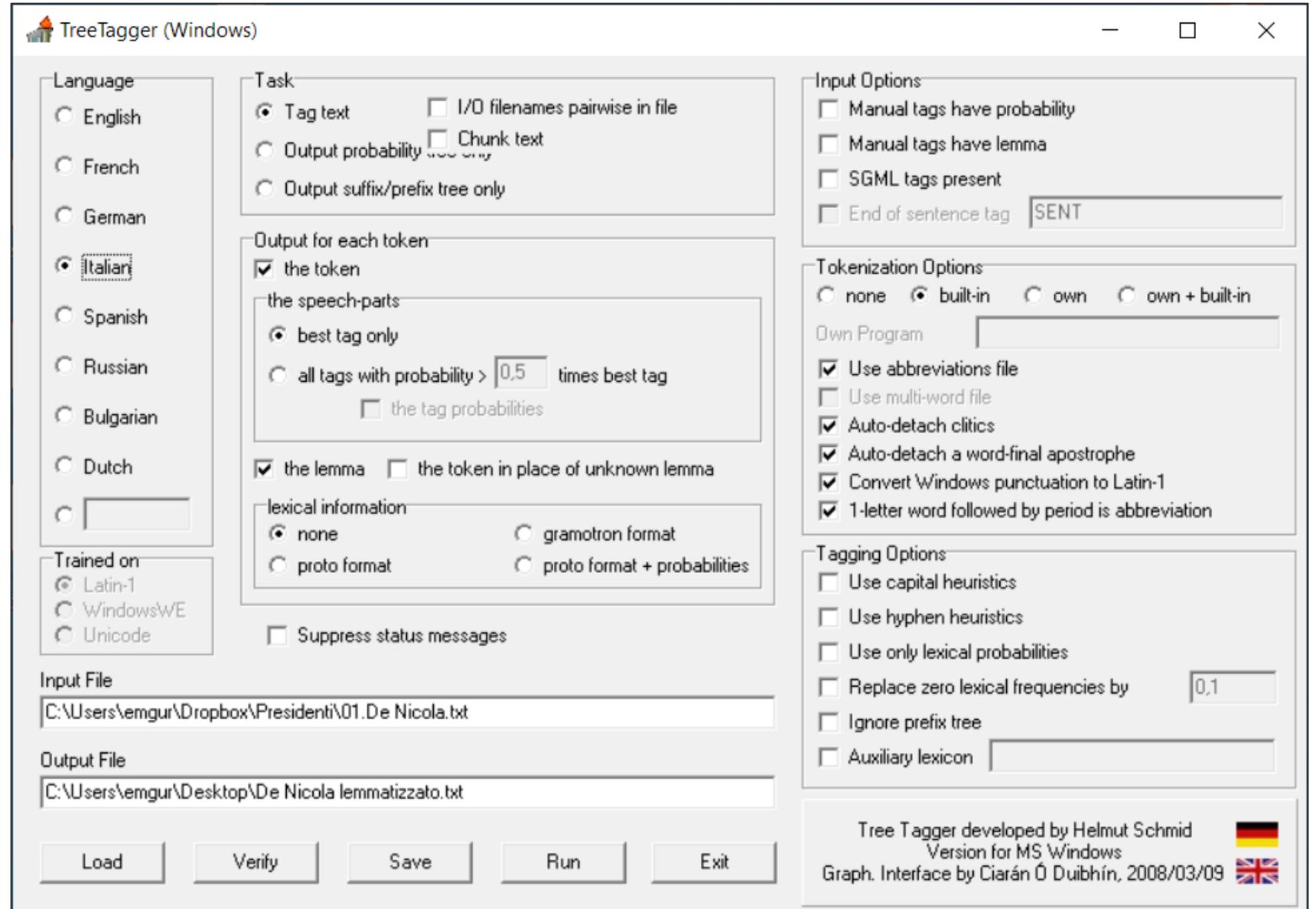
Ambiguità risolte  
considerando il contesto.  
Es. Marco porta una torta

# LEMMATIZZAZIONE

## TreeTagger

Demo:

<https://cental.uclouvain.be/treetagger/>



# LEMMATIZZAZIONE (2)

	A	B	LEMMA
1	FORMA	POS	LEMMA
2	Giuro	VER:pres	giurare
3	davanti	ADJ	davanti
4	al	PRE:det	al
5	popolo	NOM	popolo
6	italiano	ADJ	italiano
7	,	PON	,
8	per	PRE	per
9	mezzo	NOM	mezzo
10	della	PRE:det	del
11	Assemblea	NOM	assemblea
12	Costituente	ADJ	costituente
13	,	PON	,
14	che	PRO:rela	che
15	ne	ADV	ne
16	è	VER:pres	essere
17	la	DET:def	il
18	diretta	ADJ	diritto
19	e	CON	e
20	legittima	ADJ	legittimo
21	rappresentanza	NOM	rappresentanza
22	,	PON	,
23	di	PRE	di
24	compiere	VER:infi	compiere
25	il	DET:def	il
26	mia	PRO:poss	mio
27	breve	ADJ	breve
28	ma	CON	ma
29	intenso	ADJ	intenso
30	missione	NOM	missione
31	di	PRE	di
32	Capo	NOM	capo
33	provvisorio	ADJ	provvisorio
34	dello	PRE:det	del
35	Stato	NOM	stato
36	inspirandomi	VER:geru	inspirare
37	ad	PRE	ad
38	un	DET:indef	un
39	solo	ADV	solo
40	ideale	ADJ	ideale
41	:	PON	:
42	di	PRE	di
43	servire	VER:infi	servire
44	con	PRE	con
45	fedeltà	NOM	fedeltà

Italian tagset used in the TreeTagger parameter file  
(Copyright Prof. Achim Stein, University of Stuttgart)

ABR	abbreviation
ADJ	adjective
ADV	adverb
CON	conjunction
DET:def	definite article
DET:indef	indefinite article
FW	foreign word
INT	interjection
LS	list symbol
NOM	noun
NPR	name
NUM	numeral
PON	punctuation
PRE	preposition
PRE:det	preposition+article
PRO	pronoun
PRO:demo	demonstrative pronoun
PRO:indef	indefinite pronoun
PRO:inter	interrogative pronoun
PRO:pers	personal pronoun
PRO:poss	possessive pronoun
PRO:refl	reflexive pronoun
PRO:rela	relative pronoun
SENT	sentence marker
SYM	symbol
VER:cimp	verb conjunctive imperfect
VER:cond	verb conditional
VER:cpre	verb conjunctive present
VER:futu	verb future tense
VER:geru	verb gerund
VER:impe	verb imperative
VER:impf	verb imperfect
VER:infi	verb infinitive
VER:pper	verb participle perfect
VER:ppre	verb participle present
VER:pres	verb present
VER:refl:infi	verb reflexive infinitive
VER:remo	verb simple past

## Testo originale

Giuro davanti al popolo italiano, per mezzo della Assemblea Costituente, che ne è la diretta e legittima rappresentanza, di compiere la mia breve ma intensa missione di Capo provvisorio dello Stato ispirandomi ad un solo ideale: di servire con fedeltà e con lealtà il mio Paese.

Per l'Italia si inizia un nuovo periodo storico di decisiva importanza. All'opera immane di ricostruzione politica e sociale dovranno concorrere, con spirito di disciplina e di abnegazione, tutte le energie vive della nazione, non esclusi coloro i quali si siano purificati da fatali errori e da antiche colpe.

## Testo lemmatizzato

giurare davanti al popolo italiano, per mezzo del assemblea costituente, che ne essere il diretto e legittimo rappresentanza, di compiere il mio breve ma intenso missione di capo provvisorio del stato ispirare ad un solo ideale: di servire con fedeltà e con lealtà il mio paese. per il Italia si iniziare un nuovo periodo storico di decisivo importanza. al opera immane di ricostruzione politico e sociale dovere concorrere, con spirito di disciplina e di abnegazione, tutto il energia vivo del nazione, non escluso **colorare** il quale si essere purificare da fatale errore e da antico colpa

# LEMMATIZZAZIONE (3)

CONCORDANCE Italian parliamentary debates (ParlaMint 2.1) Get more space ? ! + ☆

lemma **votare** • 19,560  
636.67 per million tokens • 0.064%

Details Left context KWIC Right context

1	<input type="checkbox"/> <a href="#">CalderoliRobert...</a>	... da interpretarsi questo articolo, leggendolo al contrario, nel caso tra i nostri banchi sieda e possa venire interpretare questo articolo, leggere al contrario, nel caso tra i nostri banchi sedere e potere venire	<b>votata</b> votare	una persona che è già rappresentante del Governo e che lettura dia lei dell'articolo 13. La mia per una persona che essere già rappresentante del Governo e che lettura dare lei del articolo 13. la mia per	
2	<input type="checkbox"/> <a href="#">ColomboEmilio • ...</a>	... che il presidente provvisorio De Martino, a seguito della votazione di ballottaggio tra i due candidati più che il presidente provvisorio De Martino, a seguito della votazione di ballottaggio tra il 2 candidato più	<b>votati</b> votare	nel terzo scrutinio, ha letto il seguente risultato: "Senatori presenti: 325. Senatori votanti: 325. Sc nel terzo scrutinio, avere leggere il seguente risultato: " Senatori presente: 325. senatore votante: 325. Sc	
3	<input type="checkbox"/> <a href="#">GrassoPietro • ...</a>	... le operazioni di voto il senatore Casini chiede che gli sia data un'altra scheda poiché ha sbagliato a a operazione di voto il senatore Casini chiedere che gli essere dare una altra schedare poiché avere sbagliare a	<b>votare</b> votare	. Il Presidente, acconsentendo alla richiesta: "C' è qualche senatore che vota per la prima volta... i . Il Presidente, acconsentire alla richiesta: " c' essere qualche senatore che votare per la primo volta ... i	
4	<input type="checkbox"/> <a href="#">GrassoPietro • ...</a>	... hé ha sbagliato a votare. Il Presidente, acconsentendo alla richiesta: "C' è qualche senatore che hé avere sbagliare a votare. il Presidente, acconsentire alla richiesta: " c' essere qualche senatore che	<b>vota</b> votare	per la prima volta... in Senato certamente!"). PRESIDENTE. Dichiaro chiusa la votazione e invito i se per la primo volta ... in Senato certamente !"). presidente. dichiarare chiudere la votazione e invito il se	
5	<input type="checkbox"/> <a href="#">CompagnaLuigi • ...</a>	... a e disciplina di bilancio". Del resto, quando in quest'Aula - mi pare poco prima dell'estate - dovemmo i e disciplina di bilancio ". del resto, quando in questo Aula - mi parere poco primo del estate - dovere	<b>votare</b> votare	la fiducia sul provvedimento del ministro Passera, lo stesso Ministro ammetteva che spazi di crescita la fiducia sul provvedimento del ministro Passera, il stesso Ministro ammettere che spaziare di crescita	
6	<input type="checkbox"/> <a href="#">CompagnaLuigi • ...</a>	... io sostenuto con passione il Governo Berlusconi, con la stessa passione - lei lo ricordava - le abbiamo sostenere con passione il Governo Berlusconi, con la stesso passione - lei lo ricordare - la avere	<b>votato</b> votare	la fiducia quel 17 novembre e, molte volte, anche in occasioni successive. Detto questo, però, "I' Eu la fiducia quello 17 novembre e, molto volta, anche in occasione successivo. dire questo, però, "I' Eu	
7	<input type="checkbox"/> <a href="#">CrimiVitoClaudi...</a>	... blea del Senato di disporre l'istituzione di una Commissione speciale, quindi che sia l'Assemblea a blea del Senato di disporre il istituzione di una Commissione speciale, quindi che essere l' Assemblea a	<b>votarla</b> votare	, e di voler disporre, ai sensi dell'articolo 77 del Regolamento, con una deliberazione dell' Assemblea , e di volere disporre, al senso del articolo 77 del Regolamento, con una deliberazione dell' Assemblea	
8	<input type="checkbox"/> <a href="#">CrimiVitoClaudi...</a>	... ed è stata modificata. In merito all'urgenza, tengo a sottolineare che la risoluzione che andremo a ed essere essere modificare . in merito al urgenza, tenere a sottolineare che la risoluzione che andare a	<b>votare</b> votare	martedì è un provvedimento di una portata non indifferente, al quale non vogliamo sottrarci, anzi v martedì essere un provvedimento di una portare non indifferente, al quale non volere sottrarre, anzi v	
9	<input type="checkbox"/> <a href="#">DivinaSergio • ...</a>	... ssiamo anche dire che, se esiste un Governo tecnico, significa che qualcuno l' ha sostenuto e anche votere anche dire che, se esistere un Governo tecnico, significare che qualcuno il avere sostenere e anche	<b>votato</b> votare	. Pertanto, questo fuggi fuggi dalle responsabilità sembra dovuto al fatto che all'ultima ora tutti portanc . pertanto, questo fuggire fuggire dalla responsabilità sembrare dovere al fatto che al ultimo ora tutti portare	
10	<input type="checkbox"/> <a href="#">FerraraMario • ...</a>	... entando la legge di contabilità. Infatti oggi - lo dico a chi ha esperienza e a chi non ce l' ha - stiamo tentare la legge di contabilità . infatti oggi - lo dicare a chi avere esperienza e a chi non ce il avere - stare	<b>votando</b> votare	la possibilità che si copra con il debito pubblico, cosa specificamente negata dalla legge di contabilità, la possibilità che si coprire con il debito pubblico, cosa specificamente negare dalla legge di contabilità,	
11	<input type="checkbox"/> <a href="#">MaranAlessandro...</a>	... ontare questo importante decreto-legge. Signor Presidente, colleghi, il Gruppo Scelta Civica per l'Italia ontare questo importante decreto-legge. signore Presidente, collega, il Gruppo Scelta civico per il Italia	<b>voterà</b> votare	a favore della proposta di risoluzione n. 2, che abbiamo sottoscritto insieme alla maggioranza dei Gru a favore della proposta di risoluzione n. 2, che avere sottoscrivere insieme alla maggioranza del Gru	
12	<input type="checkbox"/> <a href="#">CappellettiEnri...</a>	... ori, confermo naturalmente che, come Movimento 5 Stelle, per coerenza con la nostra linea politica di io, confermare naturalmente che, come Movimento 5 Stelle, per coerenza con la nostra linea politico di	<b>votare</b> votare	le idee e i provvedimenti, abbiamo ritenuto opportuno ritirare la nostra proposta di risoluzione per conv le idee e il provvedimento, avere ritenere opportuno ritirare la nostra proposta di risoluzione per conv	
13	<input type="checkbox"/> <a href="#">AzzolliniAntoni...</a>	... no discusso a lungo, e naturalmente il Gruppo del Popolo della Libertà è impegnato non soltanto a i discusso a lungo, e naturalmente il Gruppo del Popolo della Libertà essere impegnare non soltanto a	<b>votare</b> votare	convintamente la proposta di risoluzione all'esame ma anche e soprattutto a verificare che l'azione pro convintamente la proposta di risoluzione al esame ma anche e soprattutto a verificare che il azione pro	
14	<input type="checkbox"/> <a href="#">SangalliGianCar...</a>	... sidente, signor Ministro, rappresentanti del Governo, colleghi senatori, il Gruppo Partito Democratico sidente, signore Ministro, rappresentante del Governo, collega senatorio, il Gruppo Partito democratico	<b>voterà</b> votare	a favore della proposta di risoluzione comune, e contemporaneamente apprezza la proposta che il Go a favore della proposta di risoluzione comune, e contemporaneamente apprezzare la proposta che il Go	
15	<input type="checkbox"/> <a href="#">CirinnaMonica • ...</a>	... per quanto riguarda la rappresentanza. Per la prima volta, tra un mese e mezzo, in questo Comune si per quanto riguardare la rappresentanza . per la primo volta, tra un mese e mezzo, in questo comune si	<b>voterà</b> votare	con un numero ridotto di consiglieri comunali, che, con i decreti, è sceso a 48. Per fortuna - lo di con un numero ridurre di consigliere comunale, che, con il decreto, essere scendere a 48. per fortuna - lo di	
16	<input type="checkbox"/> <a href="#">BisinellaPatriz...</a>	... a presenza di un rappresentante del Governo nel momento in cui si dovevano dibattere, illustrare e poi a presenza di un rappresentante del Governo nel momento in cui si dovere dibattere, illustrare e poi	<b>votare</b> votare	gli emendamenti, nonché la proposta di parere (come poi è stato fatto), non si è potuto che fa gli emendamenti, nonché la proposta di parere ( come poi essere essere fare ), non si essere potere che far	
17	<input type="checkbox"/> <a href="#">PalermoFrancesc...</a>	... enti sub-statali che, fino adesso, ha molto caratterizzato l'attuazione del federalismo fiscale. Quindi, ente sub-statali che, fino adesso, avere molto caratterizzare il attuazione del federalismo fiscale. quindi,	<b>votando</b> votare	favorevolmente, auspichiamo che in futuro le autonomie territoriali tutte possano essere valorizzate e - favorevolmente, auspicare che in futuro la autonomia territoriale tutte potere essere valorizzare e -	
18	<input type="checkbox"/> <a href="#">FerraraMario • ...</a>	... ntanza democratica, alla capacità decisionale, ai poteri del Governo e alla figura del Capo dello Stato. ntanza democratico, alla capacità decisionale, al potere del Governo e alla figura del Capo del Stato.	<b>Votiamo</b> votare	dunque a favore di tale provvedimento, ma speriamo che esso sia coniugato, nel più stretto margine dunque a favore di tale provvedimento, ma sperare che esso essere coniugare, nel più stretto margine	



# MISURE LESSICOMETRICHE ESSENZIALI

Termini		Parola	Conteggio	Relativa
<input type="checkbox"/>	1	stato	105	4,000
<input type="checkbox"/>	2	libertà	83	3,162
<input type="checkbox"/>	3	repubblica	79	3,009
<input type="checkbox"/>	4	ogni	76	2,895
<input type="checkbox"/>	5	popolo	75	2,857
<input type="checkbox"/>	6	costituzione	74	2,819
<input type="checkbox"/>	7	essere	71	2,705
<input type="checkbox"/>	8	paese	64	2,438
<input type="checkbox"/>	9	vita	61	2,324
<input type="checkbox"/>	10	politica	59	2,247
<input type="checkbox"/>	11	pace	59	2,247
<input type="checkbox"/>	12	Italia	55	2,095
<input type="checkbox"/>	13	giustizia	55	2,095
<input type="checkbox"/>	14	sempre	52	1,981
<input type="checkbox"/>	15	parlamento	52	1,981
<input type="checkbox"/>	16	democrazia	52	1,981
<input type="checkbox"/>	17	sociale	51	1,943
<input type="checkbox"/>	18	presidente	50	1,905
<input type="checkbox"/>	19	senza	46	1,752
<input type="checkbox"/>	20	forze	46	1,752
<input type="checkbox"/>	21	nazionale	45	1,714
<input type="checkbox"/>	22	lavoro	45	1,714
<input type="checkbox"/>	23	solo	44	1,676
<input type="checkbox"/>	24	italiano	43	1,638
<input type="checkbox"/>	25	mondo	42	1,600
<input type="checkbox"/>	26	deve	41	1,562
<input type="checkbox"/>	27	responsabi...	40	1,524
<input type="checkbox"/>	28	può	40	1,524
<input type="checkbox"/>	29	valori	39	1,486
<input type="checkbox"/>	30	onorevoli	39	1,486
<input type="checkbox"/>	31	saluto	38	1,448
<input type="checkbox"/>	32	parte	38	1,448
<input type="checkbox"/>	33	comunità	38	1,448

**Frequenza Assoluta:** numero di occorrenze di una parola nel corpus

**Frequenza Relativa:** rapporto tra la frequenza relativa di una parola e le dimensioni del corpus (token); può essere **normalizzata** moltiplicandola per una base prestabilita (standard: 1 milione)

↓  
Perché è importante?

**Rango** Ordine nella lista delle parole ordinate secondo la loro frequenza assoluta

Distribuzione normale:

- le parole più frequenti sono parole grammaticali – è atteso, non è un risultato dell'analisi! Possono essere rimosse?
- il numero delle occorrenze declina rapidamente, già tra rango 1 e 2
- pochi elementi con frequenza altissima, molti con frequenza bassa – specialmente **hapax**.

↓  
% hapax è un indicatore rilevante.

Dovrebbe essere <50%

# MISURE LESSICOMETRICHE ESSENZIALI (2)

**Media** (*mean*) Somma di tutti i valori di una variabile divisa per il numero totale dei valori.

**Mediana** (*median*) È il valore che divide la distribuzione a metà dopo aver ordinato i valori dal più piccolo al più grande.

**TTR** Rapporto tra type e token, è un indice di ricchezza lessicale. Esprimibile come una percentuale moltiplicandolo per 100. Risponde alla domanda: quanto è vario il lessico utilizzato in un testo? ➡ Più è alto il valore, più il lessico è ricco.

Documents					
	Title ↑	Words	Types	Ratio	Words/Sentence
1	01.De Nicola	549	342	62%	34.3
2	02.Einaudi	825	464	56%	55.0
3	03.Gronchi	2,357	1,015	43%	37.4
4	04.Segni	1,684	670	40%	31.2
5	05.Saragat	1,530	662	43%	30.0
6	06.Leone	2,094	904	43%	30.3
7	07.Pertini	1,338	615	46%	20.6
8	08.Cossiga	4,425	1,516	34%	39.9
9	09.Scalfaro	3,456	1,270	37%	27.4
10	10.Ciampi	1,920	854	44%	27.8
11	11.Napolitano	3,738	1,442	39%	34.6
12	12.Mattarella	2,336	1,029	44%	17.4

La misura ha valore quando impiegata per confrontare corpora o sub corpora.

Poiché il numero dei token influenza quello dei type è possibile ricorrere al

**TTR standardizzato** [TTR medio calcolato su segmenti di lunghezza predefinita] – ha altre debolezze...

# INDIVIDUAZIONE DEL LESSICO RILEVANTE

Come selezionare le forme su cui concentrare l'analisi?

Si possono considerare diversi soglie e pesi.

Quanti e quali dipende degli scopi della ricerca.

- Frequenza → soglia di frequenza, o un determinato tasso di copertura del corpus.

Inoltre: includere o escludere le stop words?

- Estrazione delle keyword

- Collocazioni

# ESTRAZIONE DELLE KEYWORD

Per evidenziare le peculiarità di un testo è necessario il confronto tra corpora o sub corpora.

Si possono impiegare diverse misure di keyness (= rilevanza)

## 1. Simple maths

Compara il **focus corpus** con un **corpus di riferimento**.

Considera il rapporto tra la frequenza normalizzata di una parola nel focus corpus e in quello di riferimento  $\longrightarrow$  il numero ottenuto ne rappresenta la **keyness**:  
più il valore è alto, più la parola è distintiva.

Misura spesso impiegata per estrarre la terminologia.

Raffinabile con l'aggiunta di una costante: 1 [trova parole di bassa frequenza], 100, 1000 [trova parole di alta frequenza].

$$\frac{fpm_{rmfocus} + N}{fpm_{rmref} + N}$$

# SIMPLE MATHS

## Funzione **keywords** in (No)Sketch Engine Confronto corpus ParlaMint vs ItTenTen20

SINGLE-WORDS ✓ MULTI-WORD TERMS ✓

reference corpus: Italian Web 2020 (itTenTen20) (items: 68,638)

	Lemma	Frequency per million?				Lemma	Frequency per million?		
		Focus	Reference	Score ?			Focus	Reference	Score ?
1	indico	812.28	0.06	770.5 ...	14	riformulazione	140.00	1.09	67.6 ...
2	scrutinio	1,156.82	3.37	264.7 ...	15	in-aut	59.86	0.01	60.1 ...
3	simultaneo	1,122.22	3.53	247.9 ...	16	fi-pdl	57.42	0.09	53.6 ...
4	senatrice	716.26	2.98	180.2 ...	17	prescritto	191.85	2.73	51.7 ...
5	votazione	2,781.73	14.86	175.5 ...	18	senato	1,658.77	32.19	50.0 ...
6	nominale	1,017.60	6.33	138.9 ...	19	relatrice	134.46	1.76	49.0 ...
7	senatore	3,007.66	21.07	136.3 ...	20	malan			
8	emendamento	2,952.88	21.24	132.8 ...	21	zanda			
9	petris	160.89	0.39	116.1 ...	22	airola			
10	santangelo	135.96	0.66	82.4 ...	23	caliendo			
11	capigruppo	119.03	0.58	75.8 ...	24	calderoli			
12	senatorio	190.64	1.53	75.7 ...	25	crimi			
13	decreto-legge	548.59	6.36	74.6 ...	26	pronunziare			

SINGLE-WORDS ✓ MULTI-WORD TERMS ✓

reference corpus: Italian Web 2020 (itTenTen20) (items: 1,700,169)

	Term	Frequency per million?				Term	Frequency per million?		
		Focus	Reference	Score ?			Focus	Reference	Score ?
1	scrutinio simultaneo	1,121.40	0.02	1,104.8 ...	18	ripresa della discussione del disegno	128.25	< 0.01	128.9 ...
2	votazione nominale con scrutinio simultaneo	998.79	0.01	987.7 ...	19	parere contrario	227.75	1.01	114.1 ...
3	votazione nominale con scrutinio	998.76	0.01	987.6 ...	20	rappresentante del governo	220.17	0.97	112.4 ...
4	votazione nominale	1,006.83	0.14	881.3 ...	21	parte dell' emendamento	104.39	0.03	102.4 ...
5	scrutinio simultaneo dell' emendamento	469.01	< 0.01	469.4 ...	22	signore ministro	190.09	0.99	95.8 ...
6	procedimento elettronico	426.40	0.09	391.1 ...	23	richiesta di votazione	95.40	0.02	94.8 ...
7	signora presidente	442.45	0.19	372.8 ...	24	parere contrario ai sensi dell' articolo	93.87	< 0.01	94.6 ...
8	signore presidente	1,584.07	5.25	253.5 ...	25	parere contrario ai sensi	93.87	< 0.01	94.6 ...
9	senatore segretario	198.00	< 0.01	197.7 ...	26	espresso parere contrario	102.27	0.09	94.4 ...
10	votazione dell' emendamento	204.58	0.06	194.4 ...	27	resoconto della seduta	101.39	0.12	91.8 ...
11	prescritto numero di senatori	188.40	< 0.01	187.9 ...	28	processo verbale	153.08	0.70	90.7 ...
12	prescritto numero	188.79	0.02	186.3 ...	29	seduta odierna	128.90	0.47	88.5 ...
13	numero di senatori	189.86	0.04	183.0 ...	30	b al resoconto della seduta	77.89	< 0.01	78.8 ...
14	discussione del disegno di legge	179.45	0.10	163.6 ...	31	b al resoconto	77.89	< 0.01	78.8 ...

# SIMPLE MATHS (2)

## Funzione **keywords** in (No)Sketch Engine

### Confronto sub corpora games vs news

SINGLE-WORDS ✓ MULTI-WORD TERMS ✕

reference corpus: Italian Web 2020 (itTenTen20) subcorpus: Topic news (items: 187,115)

Word	Score ?	Word	Score ?	Word	Score ?	Word	Score ?	Word	Score ?
1 pokémon	1,142.5 ...	11 poké	142.8 ...	21 lucinda	93.7 ...	31 charizard	77.8 ...	41 capcom	70.8 ...
2 ash	733.4 ...	12 playstation	130.8 ...	22 cest	92.2 ...	32 alola	76.2 ...	42 oak	70.3 ...
3 pikachu	310.7 ...	13 multiplayer	130.8 ...	23 switch	91.7 ...	33 unima	76.0 ...	43 zaffiro	70.2 ...
4 rocket	235.1 ...	14 dungeon	119.2 ...	24 ball	91.3 ...	34 mewtwo	74.0 ...	44 eevee	66.2 ...
5 gameplay	231.7 ...	15 misty	113.8 ...	25 wii	89.7 ...	35 rpg	74.0 ...	45 e3	64.8 ...
6 pokédex	167.7 ...	16 meowth	112.7 ...	26 kanto	86.6 ...	36 jessie	73.8 ...	46 ubisoft	64.0 ...
7 tapatalk	164.4 ...	17 dlc	108.7 ...	27 asd	80.3 ...	37 evil	72.5 ...	47 ps3	64.0 ...
8 xbox	161.6 ...	18 capopalestra	107.4 ...	28 xd	79.5 ...	38 johto	72.4 ...	48 goku	63.7 ...
9 nintendo	146.9 ...	19 glitch	101.7 ...	29 sinnoh	78.2 ...	39 hoenn	72.2 ...	49 dragon	62.0 ...
10 brock	145.4 ...	20 ps4	98.5 ...	30 manga	78.0 ...	40 steam	71.5 ...	50 console	61.4 ...

Rows per page: 50 1-50 of 1,000 1 / 20

Per evidenziare le peculiarità di un testo è necessario il confronto tra corpora o sub corpora.

Si possono impiegare diverse misure di keyness (= rilevanza)

## 2. TF-IDF

Mette in relazione la frequenza di una parola in un sub corpus (TF), con la frequenza inversa della parola in tutti i documenti (IDF).  $TF \times \log_{10} \frac{M}{m}$    
 M = numero di testi   
 m = testi in cui la parola è attestata



Se una parola è molto frequente in un testo, e assente in quasi tutti gli altri, allora è fortemente distintiva → lessico peculiare di un sub corpus

 Sommarario

Questo corpus contiene 12 documenti con 26,252 totale parole e con 5,290 forme di parola uniche. Creato circa un giorno fa.

Lunghezza del documento: 

- Il più lungo/a: 08.Cossiga (4425); 11.Napolitano (3738); 09.Scalfaro (3456); 03.Gronchi (2357); 12.Mattarella (2336)
- Il più corto/a: 01.De Nicola (549); 02.Einaudi (825); 07.Pertini (1338); 05.Saragat (1530); 04.Segni (1684)

Densità del vocabolario: 

- Maggiore: 01.De Nicola (0.623); 02.Einaudi (0.562); 07.Pertini (0.460); 10.Ciampi (0.445); 12.Mattarella (0.440)
- Minore: 08.Cossiga (0.343); 09.Scalfaro (0.367); 11.Napolitano (0.386); 04.Segni (0.398); 03.Gronchi (0.431)

Average Words Per Sentence: 

- Maggiore: 02.Einaudi (55.0); 08.Cossiga (39.9); 03.Gronchi (37.4); 11.Napolitano (34.6); 01.De Nicola (34.3)
- Minore: 12.Mattarella (17.4); 07.Pertini (20.6); 09.Scalfaro (27.4); 10.Ciampi (27.8); 05.Saragat (30.0)

Readability Index: 

- Maggiore: 06.Leone (17.081); 03.Gronchi (16.278); 08.Cossiga (15.722); 10.Ciampi (15.606); 11.Napolitano (15.371)
- Minore: 01.De Nicola (13.927); 02.Einaudi (13.940); 09.Scalfaro (13.977); 07.Pertini (14.228); 04.Segni (14.674)

Parole più frequenti nel corpus

- **stato** (105); **libertà** (83); **repubblica** (79); **ogni** (76); **popolo** (75)

Parole caratteristiche (in relazione al resto del corpus)

- 01.De Nicola: **ordinario** (2), **dovranno** (2), **umiliazione** (1), **tristezza** (1), **tributario** (1).
- 02.Einaudi: **pocci** (4), **opinione** (3), **aver** (3), **innanzi** (2), **gioia** (2).
- 03.Gronchi: **dell'iniziativa** (4), **trasformazione** (5), **privata** (5), **finanziaria** (3), **esistenza** (3).
- 04.Segni: **reverente** (3), **supremi** (4), **aspirazione** (4), **prestato** (2), **ordinato** (2).
- 05.Saragat: **trae** (3), **salvaguardia** (2), **inviolabile** (2), **attuazione** (2), **affidatomi** (2).
- 06.Leone: **complesso** (4), **sovrattutto** (3), **aspettative** (3), **tratta** (3), **varietà** (2).
- 07.Pertini: **violenza** (7), **difendere** (5), **bisogna** (3), **vada** (2), **tentasse** (2).
- 08.Cossiga: **speranza** (10), **perini** (3), **l'esercizio** (3), **intellettuali** (3), **generazione** (3).
- 09.Scalfaro: **vorrei** (4), **denominatore** (4), **gente** (7), **dell'uomo** (11), **occorre** (13).
- 10.Ciampi: **stabilità** (6), **europea** (9), **signor** (4), **europeo** (5), **modello** (3).
- 11.Napolitano: **ruolo** (8), **rilancio** (3), **rappresentative** (3), **raccogliere** (3), **opposti** (3).
- 12.Mattarella: **significa** (15), **ragazzi** (3), **napolitano** (3), **globali** (3), **volto** (7).

## Voyant Tools

Le **parole caratteristiche** vengono individuate sulla base del TF-IDF.

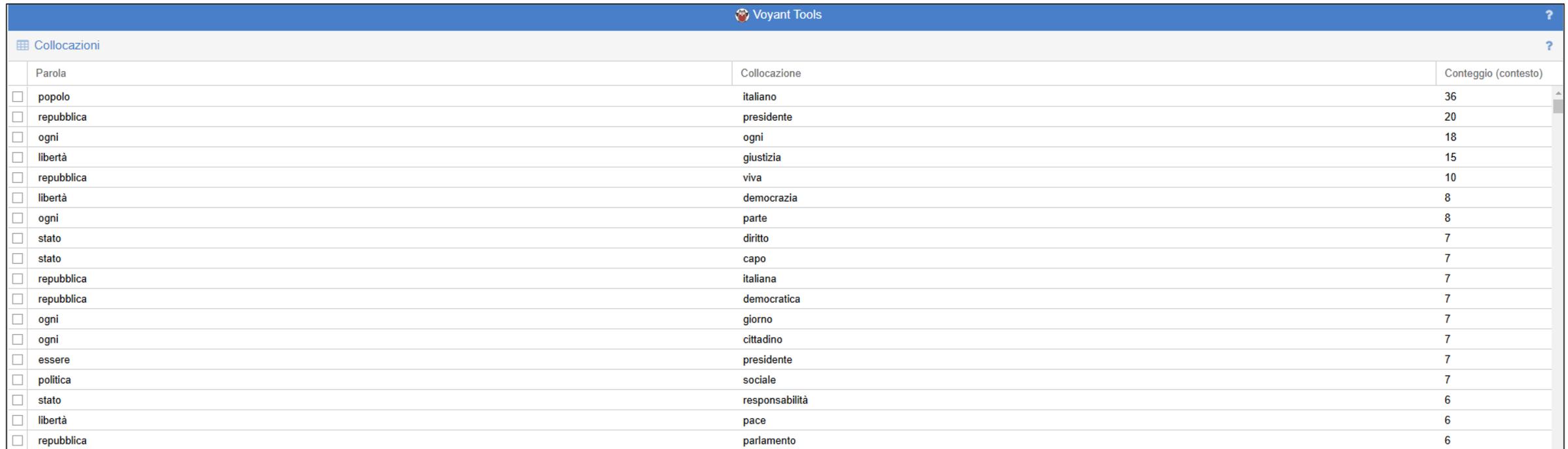
# MISURE DI COLLOCAZIONE

Individuano co-occorrenze, e ne segnalano il peso.

**1. Frequenza osservata** Numero di co-occorrenze tra due parole in contesti di lunghezza predefinita.

È necessario rimuovere le stopwords per ottenere risultati utili.

Es.: funzioni 'sintagmi' e 'collocazioni' di Voyant tools



The screenshot shows the 'Collocazioni' (Collocations) tool in Voyant Tools. It displays a table with three columns: 'Parola' (Word), 'Collocazione' (Collocation), and 'Conteggio (contesto)' (Count (context)). The table lists various words and their associated collocations with their respective counts.

Parola	Collocazione	Conteggio (contesto)
<input type="checkbox"/> popolo	italiano	36
<input type="checkbox"/> repubblica	presidente	20
<input type="checkbox"/> ogni	ogni	18
<input type="checkbox"/> libertà	giustizia	15
<input type="checkbox"/> repubblica	viva	10
<input type="checkbox"/> libertà	democrazia	8
<input type="checkbox"/> ogni	parte	8
<input type="checkbox"/> stato	diritto	7
<input type="checkbox"/> stato	capo	7
<input type="checkbox"/> repubblica	italiana	7
<input type="checkbox"/> repubblica	democratica	7
<input type="checkbox"/> ogni	giorno	7
<input type="checkbox"/> ogni	cittadino	7
<input type="checkbox"/> essere	presidente	7
<input type="checkbox"/> politica	sociale	7
<input type="checkbox"/> stato	responsabilità	6
<input type="checkbox"/> libertà	pace	6
<input type="checkbox"/> repubblica	parlamento	6

Individuano co-occorrenze, e ne segnalano il peso.

## 2. LogDice

Mette in relazione la frequenza relativa delle parole con quella della collocazione.



Assegnato un valore alto se le parole considerate occorrono (quasi) esclusivamente nella collocazione considerata → associazione forte.

Es.: funzione 'word sketch' di Sketch Engine.

# MISURE DI COLLOCAZIONE (2)

←	☰	🔍	✕	←	☰	🔍	✕	←	☰	🔍	✕	←	☰	🔍	✕	←	☰	🔍	✕
<b>verbs with "immigrato" as object</b>				<b>verbs with "immigrato" as subject</b>				<b>modifiers of "immigrato"</b>				<b>nouns modified by noun "immigrato"</b>				<b>prepositional phrases with nouns</b>			
<b>sbarcare</b>	7.1	...		<b>delinquere</b>	7.5	...		<b>clandestino</b>	10.4	...		<b>Boveri</b>	8.8	...		"immigrato" in + noun	2.8%	...	
sbarcare gli immigrati				gli immigrati delinquono				immigrati clandestini				bracciante	8.0	...		"immigrato" di + noun	1.7%	...	
<b>rimpatriare</b>	7.1	...		<b>rubare</b>	7.1	...		• concentrated in: legal ?				<b>equazione</b>	7.4	...		"immigrato" nel + noun	1.1%	...	
rimpatriare gli immigrati				immigrati ci rubano il lavoro				<b>irregolare</b>	9.8	...		l'equazione immigrato uguale				"immigrato" dal + noun	0.7%	...	
<b>espellere</b>	6.8	...		<b>sbarcare</b>	6.0	...		immigrati irregolari				<b>lavoratore</b>	7.1	...		"immigrato" del + noun	0.7%	...	
espellere gli immigrati				gli immigrati sbarcano				• concentrated in: news ?				i lavoratori immigrati				"immigrato" a + noun	0.7%	...	
<b>criminalizzare</b>	6.8	...		<b>pagare</b>	5.7	...		• concentrated in: legal ?				<b>cittadino</b>	6.1	...		"immigrato" nella + noun	0.6%	...	
criminalizzare gli immigrati				immigrati pagano				<b>extracomunitario</b>	8.9	...		i cittadini immigrati				"immigrato" al + noun	0.4%	...	
<b>regolarizzare</b>	6.7	...		<b>versare</b>	5.6	...		immigrati extracomunitari				<b>genitore</b>	5.4	...		"immigrato" senza + noun	0.4%	...	
regolarizzare gli immigrati				gli immigrati versano				<b>illegale</b>	7.9	...		da genitori immigrati				"immigrato" con + noun	0.4%	...	
<b>cacciare</b>	6.2	...		<b>commettere</b>	5.5	...		immigrati illegali				<b>alunno</b>	5.4	...		"immigrato" da + noun	0.4%	...	
cacciare gli immigrati				gli immigrati commettono				<b>regolare</b>	7.7	...		alunni immigrati				"immigrato" per + noun	0.3%	...	
<b>accogliere</b>	6.2	...		<b>togliere</b>	5.4	...		• concentrated in: legal ?				<b>maggioranza</b>	4.4	...					
accogliere gli immigrati				gli immigrati tolgono				<b>musulmano</b>	7.5	...		bambini immigrati							
<b>provenire</b>	6.2	...		<b>sottrarre</b>	5.3	...		immigrati musulmani				<b>imprenditore</b>	3.6	...					
da cui provengono gli immigrati				che gli immigrati sottraggano				<b>africano</b>	7.3	...		degli imprenditori immigrati							
<b>rinchiudere</b>	6.1	...		<b>provenire</b>	5.1	...		immigrati africani				<b>ragazzo</b>	3.5	...					
rinchiusi gli immigrati				immigrati provengono				• concentrated in: news ?				ragazzi immigrati							
<b>discriminare</b>	6.0	...		<b>stuprare</b>	5.0	...		<b>residente</b>	7.2	...		<b>studente</b>	3.2	...					
discriminare gli immigrati				Immigrato stupra				immigrati residenti				degli studenti immigrati							
<b>deportare</b>	5.8	...		<b>aggreire</b>	5.0	...		<b>marocchino</b>	7.1	...									
deportare gli immigrati				immigrato aggredisce				immigrati marocchini											
<b>soccorrere</b>	5.8	...		<b>contribuire</b>	5.0	...		<b>senegalese</b>	6.9	...									
soccorrere gli immigrati				gli immigrati contribuiscono				immigrati senegalesi											
								<b>ebreo</b>	6.9	...									
								di immigrati ebrei											
								• concentrated in: reference ?											
								<b>proveniente</b>	6.8	...									
								immigrati provenienti da											
								• concentrated in: arts ?											
								• concentrated in: reference ?											

È possibile prendere in considerazione diverse misure di similarità e distanza. Le più 'semplici':

- 1. Indice di Jaccard** Considera il numero di parole comuni in rapporto al numero di parole totali.  
Il valore è compreso tra 0 (niente in comune) e 1 (identici).
- 2. Distanza di Jaccard** Esprime la diversità tra due corpora. Si calcola invertendo la misura precedente:  $1 - \text{Indice di Jaccard}$ .

# Analisi delle corrispondenze: periodizzazione dell'AGI basata sui profili lessicali delle annate

