

# **Sistemi per il recupero delle informazioni**

**Gabriele Pozzani**

**A.A. 2014/2015**

**Corso di Laurea Magistrale in  
Editoria e Giornalismo**

**Esempio completo di  
analisi di testi, indicizzazione,  
ricerca e ranking**

# Cosa vedremo

- Un esempio completo di
  - Indicizzazione
    - Costruzione del vocabolario di indicizzazione
    - Indicizzazione della collezione di documenti
  - Ricerca
    - Ranking dei risultati
      - Con modello booleano esteso

3

## La collezione di documenti

- D1:  
Pirro vince! È vera vittoria? Non si sa, ma intanto ha vinto.
- D2:  
Pirro stupisce tutti e vince. Torna **a** pace?
- D3:  
I romani sconfitti da Pirro
- D4:  
Pirro nasce nel 318a.C. e muore ad Argo nel 272a.C.

4

# Costruzione del vocabolario

## 1. divisione in parole

- La lettura dei documenti porta ad ottenere le seguenti parole

272	la	sconfitti
318	ma	si
a.C.	muore	stupisce
ad	nasce	Torna
Argo	nel	tutti
da	Non	una
e	pace	vince
grande	Pirro	vinto
ha	romani	vittoria
i	sa	

5

## Costruzione del vocabolario 2. riduzione al maiuscolo

272	la	sconfitti
318	ma	si
a.C.	muore	stupisce
ad	nasce	Torna
Argo	nel	tutti
da	non	una
e	pace	vince
grande	pirro	vinto
ha	romani	vittoria
i	sa	

6

# Costruzione del vocabolario

## 3. applicazione della stop-list

272	la	ə	sconfitti
318	ma	mɑ	si
a.C.	muore		stupisce
ad	nasce		torna
argo	nei	neɪ	tutti
da	non	nən	ma
e	pace		vince
grande	pirro		vinto
ha	romani		vittoria
i	sa		

7

# Costruzione del vocabolario

## 4. gestione accenti

- Nulla da fare

8

# Costruzione del vocabolario

## 5. stemming

272	pirro	vinto	stup
318	romani	vittoria	torn
a.c.	sa	grand	tutt
argo	sconfitti	muor	vinc
grande	si	nasc	vint
ha	stupisce	pac	vittor
muore	torna	pirr	
nasce	tutti	roman	
pace	vince	sconfitt	

Snowball stemmer, by Dr Martin Porter. <http://snowball.tartarus.org/>

# Costruzione del vocabolario

## 6. lemmatizzazione

272	romani	avere
318	sa	morire
a.c.	sconfitti	nascere
argo	si	romano
grande	stupisce	sapere
ha	torna	sconfitto
muore	tutti	stupire
nasce	vince	tornare
pace	vinto	tutto
pirro	vittoria	vincere

TreeTagger, by Helmut Schmid. <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>

# Vocabolario di indicizzazione

272	sa	nasc	avere
318	sconfitti	pac	morire
a.c.	si	pirr	nascere
argo	stupisce	roman	romano
grande	torna	sconfitt	sapere
ha	tutti	stup	sconfitto
muore	vince	torn	stupire
nasce	vinto	tutt	tornare
pace	vittoria	vinc	tutto
pirro	grand	vint	vincere
romani	muor	vittor	

11

# Indicizzazione

- Una volta costruito il vocabolario si passa alla vera e propria indicizzazione
  - Viene costruito un indice inverso
  - I termini vengono associati ai documenti in cui appaiono

12

Indice inverso (I)

13

Indice inverso (III)

14

## Vettori TF-IDF (I)

2	3	a	a	a	g	g	h	m	m	m	n	n	p	p	p	r	r	r
7	1	.	r	v	r	r	a	o	u	u	a	a	i	i	o	o	o	o
2	8	c	g	e	a	a	r	o	o	s	s	c	r	r	m	m	m	m
.	.	o	r	n	n	i	r	r	c	c	e	r	r	a	a	a	a	a
.	.	e	d	d	r	e	e	e	e	e	e	o	o	n	n	n	n	o

<b>D1</b>	0	0	0	0	1	1	0	0	0	0	0	0	0	0,25	0,25	0	0	0	
<b>D2</b>	0	0	0	0	0	0	0	0	0	0	0	0	1	1	0,25	0,25	0	0	0
<b>D3</b>	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0,25	0,25	1	1	1
<b>D4</b>	1	1	2	1	0	0	0	1	1	1	1	0	0	0	0,25	0,25	0	0	0

15

## Vettori TF-IDF (II)

s	s	s	s	s	s	s	t	t	t	t	t	y	y	y	y	y	y	y	
a	c	c	c	i	t	t	o	o	u	u	u	i	i	i	i	i	i	i	
p	o	o	o	u	u	u	r	r	t	t	t	n	n	n	n	n	n	t	
e	n	n	n	p	p	p	n	n	t	t	t	c	c	c	c	c	c	t	
r	f	f	f	i	i	i	a	a	i	o	o	e	e	e	e	e	e	o	
e	i	i	i	r	s	r	e	e	t	t	t	r	r	r	r	r	r	r	
t	t	t	t	t	t	t	t	t	t	t	t	i	i	i	i	i	i	a	
i	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o	o

<b>D1</b>	1	1	0	0	0	1	0	0	0	0	0	0	0	0,5	0,5	1	1	1
<b>D2</b>	0	0	0	0	0	1	1	1	1	1	1	1	1	0,5	0,5	0	0	0
<b>D3</b>	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<b>D4</b>	0	0	0	0	0	0	0	1	1	1	1	1	0	0	0	0	0	0

16

## Possible thesaurus (I)

- Come detto, non esiste un unico thesaurus e non esiste una regola precisa per costruirne uno
- Nel nostro esempio consideriamo di avere il seguente piccolo thesaurus

17

## Possible thesaurus (II)

- pace  
UF guerra
- vittoria  
UF sconfitta
- morire  
UF nascere
- sconfitto  
UF vinto
- roma  
RT romano
- vittoria  
RT vincere
- vittoria  
RT sconfitto
- città  
NT1 roma  
NT1 argo

18

# Esempi di ricerche

## con modello booleano esteso

- L'effettivo algoritmo (processo) per rispondere ad una ricerca dipende dal software/applicazione/motore di ricerca
  - Qui ne vedremo solo uno plausibile
- I termini della ricerca subiscono lo stesso “trattamento” delle parole nei documenti
  - Stop-list, stemming, lemmatizzazione
  - Così ANCHE i termini “correlati” a quelli inseriti vengono usati per cercare i documenti nell'indice inverso
- Useremo un modello booleano esteso con sottinteso l'operatore AND
  - Come fa Google

19

## vittoria di Pirro (I)

- Riduzione al minuscolo e stop-list: vittoria pirro
  - vittoria
    - È contenuto solo in D1 (con peso 1)
    - La sua radice e lemma sono ancora contenuti solo in D1, quindi sono “inutili”
  - Nel thesaurus è correlato a
    - vincere (RT): è contenuto anche in D2 (con peso 0,5)
    - sconfitta (UF): non compare in nessun documento
    - sconfitto (RT): compare in D3 (con peso 1)

20

## vittoria di Pirro (III)

- pirro
  - È contenuto in D1, D2, D3, D4 (sempre con peso 0,25)
  - La sua radice e lemma non portano ad altri documenti
  - Nel thesaurus non è correlato a nessun altro termine

21

## Ranking

- Ricapitolando
  - D1 ha peso totale  
 $1 + 0,25 = 1,25$
  - D2 ha peso totale  
 $0,75 \times 0,5 + 0,25 = 0,625$
  - D3 ha peso totale  
 $0,75 \times 1 + 0,25 = 1$
  - D4 ha peso totale  
0,25
- Si ottiene il ranking finale:  
D1, D3, D2, D4

22

## Risultato della ricerca: vittoria di Pirro

- (D1) Pirro vince! È vera vittoria? Non si sa, ma intanto ha vinto.
- (D3) I romani sconfitti da Pirro
- (D2) Pirro stupisce tutti e vince. Torna la pace?
- (D4) Pirro nasce nel 318a.C. e muore ad Argo nel 272a.C.