

Esercizio 1

D1: You say goodbye, I say hello

D2: You say stop, I say go

D3: Hello, hello, you say goodbye

D4: I say high, you say low

Costruire l'inverted index di questi documenti

Scegliete le stop words

Per la Query: "you goodbye" descrivete il processo di ricerca dell'inverted index e l'ordinamento secondo il modello vettoriale

Esercizio 2

QUERY: **nuovi farmaci epilessia --> nuov farmac epile**

	testo	Manual rank	Norm. Word frequency rank	Rank booleano
D1	La terapia <b>farmacologica</b> dell' <b>epilessia</b> può essere divisa, dal punto di vista storico, in due fasi, ovvero quella dei "vecchi" (dal fenobarbital - 1912, all'acido valproico - 1964), e quella dei " <b>nuovi</b> " <b>farmaci</b> (dal vigabatrim degli anni '80 ai più recenti). Questi ultimi, tuttavia vanno ulteriormente suddivisi in " <b>nuovi</b> " e " <b>nuovissimi</b> ", in quanto negli ultimi anni sono apparsi sul mercato diversi <b>farmaci</b> con caratteristiche di meccanismo d'azione ed efficacia molto interessanti.			
D2	Dall'inizio degli anni '90 sono			

	<p>stati commercializzati quasi una decina di <b>nuovi farmaci</b> antiepilettici, dopo quasi 20 anni di silenzio in questo campo. I <b>nuovi farmaci</b> hanno dimostrato globalmente efficacia almeno pari a quella dei vecchi <b>farmaci</b>, ma sicuramente una minore incidenza di effetti collaterali e quindi maggiore tollerabilità da parte del paziente. Esistono poi alcune forme particolari di <b>epilessia</b> come gli Spasmi Infantili e la Sindrome di Lennox-Gastaut in cui con i <b>nuovi farmaci</b> si ottengono risultati sicuramente migliori che in passato.</p>			
D3	<p>Si può curare l'<b>epilessia</b>?  Certamente sì, anche se solo nel 60-70% dei casi. Si usano infatti <b>farmaci</b> che controllano e bloccano la tendenza delle cellule cerebrali a produrre scariche <b>epilettiche</b>.</p>			
D4	<p>I classici <b>farmaci</b> antiepilettici sono valproato e carbamazepina (che sono spesso i farmaci di prima scelta), fenitoina e fenobarbital. Da pochi anni è disponibile una serie di <b>farmaci</b> di <b>nuova</b> generazione (felbamato, gabapentin, lamotrigina, levetiracetam, oxcarbazepina, tiagabina, topiramato, vigabatrin) usati per indicazioni particolari o per aumentare l'efficacia della terapia quando la monoterapia</p>			

	con un <b>antiepilettico</b> classico non riesce a sopprimere le crisi <b>epilettiche</b> .			
D5	I nuovi progressi nella ricerca aprono scenari interessanti per tutte le persone con <b>epilessia</b> - ha precisato il prof. Emilio Perucca, Professore Ordinario di Farmacologia, Facoltà di Medicina e Chirurgia			
D6	All'Istituto di ricerche farmacologiche Mario Negri di Milano, è attualmente in studio una <b>nuova</b> terapia genica contro <b>l'epilessia</b> .			
D7	<b>L'epilessia</b> ha una prevalenza di quasi il 10 per 1.000: in Italia si contano circa 500.000 soggetti epilettici ed ogni anno si verificano circa 25.000 <b>nuovi</b> casi.			
D8	L'incidenza è di 50 <b>nuovi</b> casi per ogni 100.000 persone ogni anno. La terapia <b>farmacologica</b> dell' <b>epilessia</b> viene sempre iniziata con un solo <b>farmaco</b> ..			
D9	<b>Epilessia</b> sotto controllo. Molti malati non prendono i <b>farmaci</b> nel modo giusto, con il rischio di avere così <b>nuovi</b> attacchi del disturbo.			
D10	Ricercatori dell'IWK Health Centre di Halifax, in Canada, hanno valutato l'incidenza di <b>epilessia</b> intrattabile in seguito ad interruzione del trattamento con farmaci antiepilettici in bambini			

epilettici della Nuova Scozia.			
--------------------------------	--	--	--

- 1) confrontate rank manuale, rank basato su normalized word frequency, e rank booleano. Elencate le vostre idee per fare di meglio, sulla base della vostra interpretazione delle esigenze di chi fa la query.
- 2) costruite l'inverted index di questi documenti per il seguente vocabolario: *bloccare, caso, cellula, crisi, cura, efficacia, epilessia, farmaco, farmacologia, epilessia, epilettico, malato, malattia, studio, terapia, trattamento, risultato*
- 3) costruite un "trie" del vocabolario

## Soluzioni

### Es. 1

	DocFr		Doc #	Freq.
go	1	→	2	1
goodbye	2	→	1	1
hello	2	→	3	1
high	1	→	1	1
i	3	→	3	2
low	1	→	4	1
say	4	→	1	1
stop	1	→	2	1
you	4	→	4	1
		→	4	1
		→	1	2
		→	2	2
		→	3	1
		→	4	2
		→	2	1
		→	1	1
		→	2	1
		→	3	1
		→	4	1

<i>tf</i>	<b>D1</b>	<b>D2</b>	<b>D3</b>	<b>D4</b>
you	1	1	1	1
goodbye	1	0	1	0

## Esercizio 2

Bloccare	1	D3,1
Caso	3	D3,1-D7,1-D8-1
Cellula	1	D3-1
Crisi	1	D4-1
Cura	1	D3-1
Efficacia	3	D1,1-D2,1-D4,1
Epilessia	9	D1,1-D2,1-D3,1-D5,1-D6,1-D7,1-D8,1-D9,1-D10,1
Epilettico	5	D2,1-D3,1-D4,3-D7,1-D10,3
Farmaco	6	D1,2-D2,4-D3,1-D4,3-D8,1-D9,1
Farmacologia	3	D1,1-D5,1-D8,1
Malato	1	D9,1
Malattia	0	
Studio	1	D6,1
Terapia	4	D1,1-D4,1-D6,1-D8,1-D10,1
Trattamento	1	D10,1
Risultato	1	D2,1

Vettori dei documenti e della query, e calcolo del rank usando  $w_{ij} = \frac{freq(w_i \text{ in } D_j)}{\max_k(freq(w_k \text{ in } D_j))}$  (senza idf)

	D1	D2	D3	D4	D5	D6	D7	D8	D9	D10	Query
Bloccare			1								
Caso			1				1	1			
Cellula			1								
Crisi				1							
Cura			1								
Efficacia	1	1		1							
Epilessia	1	1	1		1	1	1	1	1	1	
Epilettico		<b>0,33</b>	<b>0,33</b>	<b>1</b>			<b>0,33</b>		<b>0,33</b>	<b>1</b>	<b>1</b>
Farmaco	<b>0,5</b>	<b>1</b>		<b>0,75</b>				<b>1</b>	<b>1</b>		<b>1</b>
Farmacologia											
	1				1			1			
Malato									1		
Malattia											
Studio						1					
Terapia	1			1		1		1		1	
Trattamento										1	
Risultato		1									
$ D_i ^2$	4,25	4,1	5,1	4,56	2	3	2,1	5	3,1	4	2
$ D_i $	2,06	2,02	2,258	2,13	1,41	1,73	1,44	2,23	1,76	2	1,41
$D_i \cdot Q$	<b>0,5</b>	<b>1,33</b>	<b>0,33</b>	<b>1,75</b>	<b>0</b>	<b>0</b>	<b>0,33</b>	<b>1</b>	<b>1,33</b>	<b>1</b>	
$ D_i  \cdot  Q $	<b>2,91</b>	<b>4,33</b>	<b>4,83</b>	<b>4,56</b>	<b>3,02</b>	<b>3,70</b>	<b>3,10</b>	<b>4,78</b>	<b>3,76</b>	<b>4,28</b>	
COS-SIM	<b>0,17</b>	<b>0,30</b>	<b>0,068</b>	<b>0,38</b>	<b>0</b>	<b>0</b>	<b>0,10</b>	<b>0,20</b>	<b>0,35</b>	<b>0,23</b>	

(Nota: potrebbe esserci qualche errore di calcolo..)

Formula: 
$$\frac{w_{ij} \cdot w_{iq}}{\sqrt{\sum w_{ij}^2} \cdot \sqrt{\sum w_{iq}^2}}$$