INFORMATION RETRIEVAL

Week 9 – Vector Space Model

02.05.2025

Severin Mills

Today

1

Exercise Recap

- Discussion
- Questions

² Theory

- Recap: Ranked Retrieval
- Vector Space Model
- SMART Notation

Kahoot

3

• BONUS: Vector Space Model

Exercise 6: Heaps' & Zipf's Law

Recap

Heaps' Law: #terms = $k \times$ #tokens Zipf's Law: Frequency = $\frac{k}{Rank}$

Previously

- We want to rank documents based on a score
- Measure similarity between documents

Scoring

Zones in terms

Zones in postings

ETH.body	6	1	2	
ETH.title	5	3	4	
computer.body	5	1	2	
computer.title	5	1	3	
albert.title	4	2	3	
albert.author	6	1	2	
albert.body	→ 4	3	5	





Scoring



t	tf _{t,d}
foo	3
bar	2
foobar	2

Scoring

If we want to measure rarity, we use the inverse document frequency



Putting it all together

tf-idf

tf	A	В	tf-idf	Α	В
foo	5	1	foo	25	5
bar	0	4	bar	0	40
foobar	2	1	foobar	6	3



Now

Main idea: represent documents as vectors



Vectorization

Boolean vector

• Not very usable for ranking, why?



Vectorization

Use reals instead. Could be term frequency, tf-idf (preferably), etc.



Vectorization

A document is then a vector in the first quadrant, $D \in \mathbb{R}^{M}$





Vectorization

Normalize vectors



Inner Product (Renormalized)

$$\cos \theta = \frac{\vec{x}}{\|\vec{x}\|} \times \frac{\vec{y}}{\|\vec{y}\|}$$

• Good quantification of similarity



Queries as vectors

- Build vector
- Calculate cosine similarity of documents and return the nearest ones
- What about a large collection of documents?
- Doesn't matter, since query size is constant

Inverted index





SMART Notation

Term frequency		Document frequency		
n (natural)	tf _{t,d}	n (no)	1	
l (logarithm)	$1 + \log(\mathrm{tf}_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	
a (augmented)	$0.5 + \frac{0.5 \times \text{tf}_{t,d}}{\max_t(\text{tf}_{t,d})}$	p (prob idf)	$\max\{0, \log \frac{N-df_t}{df_t}\}$	
b (boolean)	$\begin{cases} 1 & \text{if } \text{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			
L (log ave)	$\frac{1 + \log(\mathrm{tf}_{t,d})}{1 + \log(\mathrm{ave}_{t \in d}(\mathrm{tf}_{t,d}))}$			



Exercise 7: Vector Space Model

BONUS TIME

- Same as previous
- Starts on 02.05 at 15:00, ends on 09.05 at 15:00
- If you already passed 2 / 3 quizzes, you got the bonus (you should get a mail soon)

Kahoot

https://create.kahoot.it/details/

<u>duplicate-of-information-</u>

retrieval-ex-07-vector-space-

models-mschoeb/ef383953-

<u>b43a-4abd-af2a-d9ebf2ad1019</u>

