

Text-based ranking

First generation, based on the text of the document and of the query

Documents as binary arrays

Until now, we represented a document as a binary array where each position corresponds to a term and each value is either zero (term is not in document) or one (term is in document).

We would score a document by overlapping it with the query

$$\text{score}(d, q) = |q \cap d|$$

It is very costly and is not very useful in term of similarity because it doesn't take into account the lengths, which usually are very different.

Two approaches for normalizing the sizes:

Signed coeff $\rightarrow \text{sim}(d, q) = \frac{|d \cap q|}{|d \cup q|}$ RESPECT TRIANGULAR INEQUALITY

Dice coeff $\rightarrow \text{sim}(d, q) = \frac{2|d \cap q|}{(|d| + |q|)}$ DOESN'T RESPECT TRIANGULAR EQUALITY

IT'S STILL NOT A GOOD WAY OF RANKING, BECAUSE WE ARE NOT CONSIDERING THE IMPORTANCE OF EACH TERM

Documents as vectors

As already said, with binary vectors we can't consider the term frequency of a document: ideally, we would like to show first the documents that tell the most about our query terms, ignoring those terms that don't add relevance (stop words).

WE WOULD LIKE TO WORK IN A VECTOR SPACE, WHERE THE AXES ARE THE TERMS AND THE POINTS OR VECTORS ARE THE DOCUMENTS AND THE QUERIES

Therefore

- 1) We want to transform queries and documents in real valued vectors and represent them in a space
- 2) We want to rank the documents according to their proximity to the query in that space

proximity = similarity of vectors \approx inverse of distance

Note = A vector space may have millions of dimensions
Note = In our case, the vectors are very sparse (most of the entries are zero)

How to estimate the proximity of two vectors?

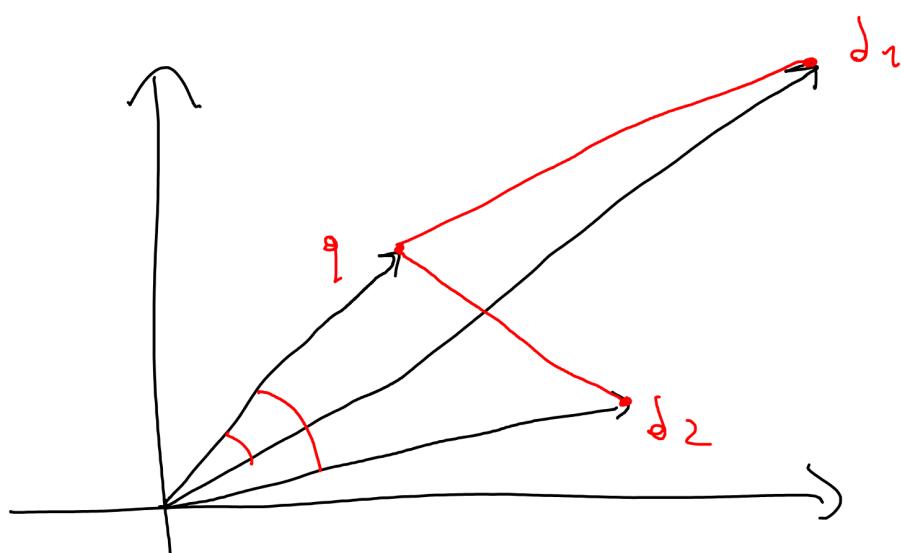
To rank a document, we need a way to compare two vectors in the vector space (namely, the query and the document).

A good proximity measure is the angle formed by the two vectors, not the distance

between them (i.e., euclidean distance)

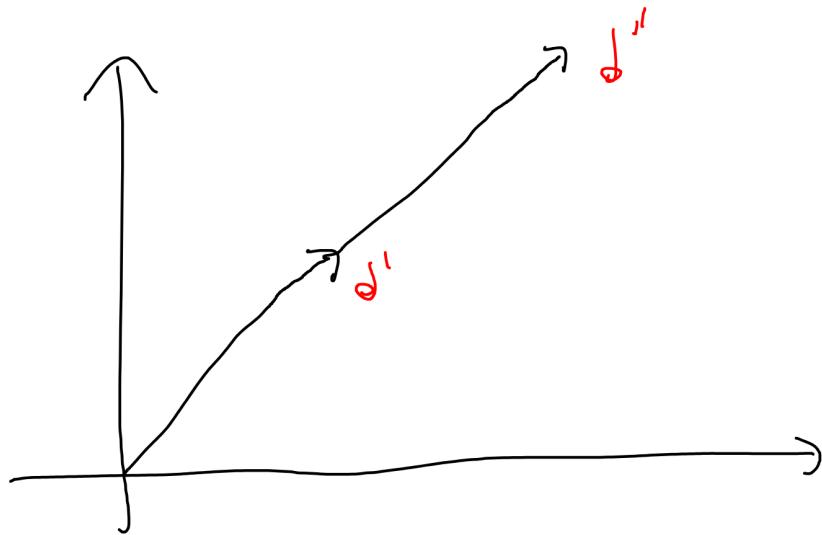
The euclidean distance may be large for vectors of different lengths, even if they are close in space.

EXAMPLE



Using euclidean distance, q is more similar to d_2 than d_1 , which is not true.

To better prove that the angle is better than the distance, here is another example



Semantically, d' and d'' have the same content:
indeed the angle is zero, but the euclidean
distance is large.

Actually, it is easier to use the cosine rather
than the angle, but the two notions are equivalent:

• USING ANGLE = VALUES CLOSER TO ZERO ARE BETTER

• USING COSINE = THE COSINE FUNCTION DECREASES MONOTONICALLY
IN THE INTERVAL $(0^\circ, 180^\circ]$ AND FURTHER VALUES ARE

BETTER $\Rightarrow \cos(0^\circ) = 1$
 $\cos(180^\circ) = -1$

But how do I get the cosine of the angle between two vectors?

COSINE SIMILARITY

$$\begin{aligned}
 \cos(\vec{q}, \vec{d}) &= \\
 &= \text{Cosine similarity of } \vec{q} \text{ & } \vec{d} = \\
 &= \text{Cosine of the angle between } \vec{q} \text{ & } \vec{d} = \\
 &= \frac{\vec{q} \cdot \vec{d}}{\|\vec{q}\| \cdot \|\vec{d}\|} = \frac{\vec{q}}{\|\vec{q}\|} \cdot \frac{\vec{d}}{\|\vec{d}\|} = \\
 &\quad \xrightarrow{\text{dot product}} \quad \xrightarrow{\text{unit vector}} \\
 &= \frac{\sum_{i=1}^{|D|} q_i d_i}{\sqrt{\sum_{i=1}^{|D|} q_i^2} \cdot \sqrt{\sum_{i=1}^{|D|} d_i^2}}
 \end{aligned}$$

NOTE = We transform the vectors into unit vectors by dividing them by their L2 norm ($\|\vec{x}\| = \sqrt{\sum_i x_i^2}$). This way the query and the documents have comparable values (weights), even if they have different lengths.

NOTE = If the two vectors are already length-normalized

$$\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d} = \sum_{i=1}^{|D|} q_i d_i$$

What are the values q_i and d_i , then?

WE CALL THEM TF-IDF WEIGHTS

TF-IDF

$$w_{t,d} = \text{weighted score of document } d \text{ for term } t = tf_{t,d} \times idf_t$$

if t is not in d or it is not relevant (i.e., t appears in all the documents) $\Rightarrow w_{t,d} = 0$

- $tf_{t,d}$ = term frequency of term t in document d =
= occurrences of t in d

- idf_t = inverted document frequency of term t =
 $= \log \left(\frac{n}{n_f} \right)$
documents in the collection
documents in the collection that contains term t
You can call it df_t for consistency

Note that we could rank a document just by the term frequency of t , but then we would be misled by stopwords terms (which appear everywhere), and that's why we also consider the inverted document frequency of t (it lowers the score if t is in many documents)

This way we replaced the binary representation of a document with a real valued vector and introduce the VECTOR SPACE MODEL (in contrast with the previous one, the Boolean model)

Example

DOCUMENT VECTOR WITH TF-IDF WEIGHTS

TERM	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	13,1	11,4	0,0	0,0	0,0	0,0
Brutus	3,0	8,3	0,0	1,0	0,0	0,0
Caesar	2,3	2,3	0,0	0,5	0,3	0,3
Calpurnia	0,0	11,2	0,0	0,0	0,0	0,0
Cleopatra	17,7	0,0	0,0	0,0	0,0	0,0
mercy	0,5	0,0	0,7	0,9	0,9	0,3
worser	1,2	0,0	0,6	0,6	0,6	0,0

Regarding the storage usage of such vectors

of course we don't store a matrix like in the example above, but we use the classic inverted list approach, with

a couple of additions:

- For every Term, the dictionary also stores the length of the corresponding posting list, so that the IDF is computed in constant time
- For every docID in the posting list, we also append the corresponding $tf_{t,d}$, which is small
- You could also add other useful metadata, or seen for the phrase query processing solutions

Okay, but how do we do the actual scoring of the documents for a given query?

COSINE_SCORE : $q \rightarrow$ Top K documents
 $\forall d \in \text{Docs}$, Returns the K best $\cos(\vec{q}, \vec{d})$ normalized

COSINE_SCORE(q) :

float scores[N] = 0 $\quad \text{\textcolor{green}{N = \# docs}}$
float lengths[N] = { $i : i = \text{size}(d) \quad d \in \text{docs}$ }
 $\quad \quad \quad \text{\textcolor{green}{\# length of the docs}}$

for t in q :

list = get posting list of t =
 $= (\text{docID}, \text{tf}_{t,d}, \text{docID})$

compute $w_{t,q}$

Dot Product PHASE [for $(d, \text{tf}_{t,d})$ in list :
compute $w_{t,d}$
 $\text{scores}[d] += w_{t,d} \cdot w_{t,q}$]

Normalization PHASE [for $(\text{score}, \text{length})$ in zip($\text{scores}, \text{lengths}$) :
 $\text{score} = \frac{\text{score}}{\text{length}}$]

return { $d : \text{scores}[d]$ is a Top K doc}

EXERCISE

Given the four texts : 1) "la belle case",
 2) "le belle case", 3) "cas la bella bella", 4) "la
 casa le case bella bella"

- compute the inverted index for the texts,
 compress with gap encoding, then with γ-code

	NORMAL	GAP	NORMAL	GAP
BELLA → 1	1, 3, 4	1, 2, 1	LA → 1, 4	1, 3
BELLE → 2		2	LE → 2, 4	2, 2
CASA → 1, 3, 4		1, 2, 1	ROSA → 4	4
CASE → 2, 4		2, 2		

$$\gamma(1) = 1, \gamma(2) = 0.1, \gamma(3) = 0.11, \gamma(4) = 0.100$$

- compute the tf-idf vectors

→ just to recall it somewhere ↓

TF-IDF	1)	2)	3)	4)	idf ←
BELLA	$1 \cdot \log\left(\frac{4}{3}\right)$	0	$2 \cdot \log\left(\frac{4}{3}\right)$	$1 \cdot \log\left(\frac{4}{3}\right)$	$\log\left(\frac{4}{3}\right)$
BELLE	0	1.2	0	0	$\log\left(\frac{4}{7}\right) \approx 2$
CASA	$1 \cdot \log\left(\frac{4}{3}\right)$	0	$1 \cdot \log\left(\frac{4}{3}\right)$	$1 \cdot \log\left(\frac{4}{3}\right)$	$\log\left(\frac{4}{3}\right)$
CASE	0	1.1	0	1.1	$\log\left(\frac{4}{2}\right) \approx 1$
LA	1.1	0	0	1.1	$\log\left(\frac{4}{2}\right) \approx 1$
LE	0	1.1	0	1.1	$\log\left(\frac{4}{2}\right) \approx 1$
ROSA	0	0	0	1.2	$\log\left(\frac{4}{1}\right) \approx 2$

- Find the text most similar to f) using cosine without normalization (simple dot product).

$$\rightarrow \text{sim}(1, f) = \log\left(\frac{f}{3}\right) \cdot \log\left(\frac{f}{3}\right) + 0 \cdot 0 + \log\left(\frac{f}{3}\right) \cdot \log\left(\frac{f}{3}\right) + \\ + 0 \cdot 1 + 1 \cdot 1 + 0 \cdot 1 + 0 \cdot 2 = \\ = 2 \cdot \log\left(\frac{f}{3}\right)^2 + 1 \approx 1.3$$

$$\rightarrow \text{sim}(2, f) = 0 \cdot \log\left(\frac{f}{3}\right) + 2 \cdot 0 + 0 \cdot \log\left(\frac{f}{3}\right) + \\ + 1 \cdot 1 + 0 \cdot 1 + 1 \cdot 1 + 0 \cdot 2 = \\ = 2$$

$$\rightarrow \text{sim}(3, f) = 2 \log\left(\frac{f}{3}\right) \cdot \log\left(\frac{f}{3}\right) + 0 \cdot 0 + \\ + \log\left(\frac{f}{3}\right) \cdot \log\left(\frac{f}{3}\right) + 0 \cdot 1 + \\ + 0 \cdot 1 + 0 \cdot 1 + 0 \cdot 2 = \\ = 3 \log\left(\frac{f}{3}\right)^2 \approx 0.52$$

2) is the most similar to f)

CONCLUSIONS ON VECTOR SPACES

The vector space is good for bag-of-words queries, but not good with operators (Boolean, wildcard, positional, proximity).

The real problem is that it is really easy for spammers to fuck up everything (just fill the document with the same terms, the Tf increases but the idf is still the same) and that's why text-based approaches are not used in modern web search engines.

ANYWAY, WE HAVEN'T FINISHED YET:



Computing the cosine is expensive due to the high-dimensionality of a vector space, but it's fast to find top scoring documents;

Solutions?

Approximate Retrieval

(Fast Top-K retrieval)

Find a set A of candidates with $K < |A| \ll N$,
compute the score for each doc (doesn't have to be exact)
and returns the K top scoring docs in A .

Set A doesn't necessarily contain all top- K docs,
but has many docs belonging to that (that's why it's
an approximation)

How do we find A ?

Consider only those docs containing at least one
query term (of course), and use one of the
following approaches:

- Only consider docs containing many of the query terms

For multi-term queries, select and score only
those docs containing most of the query terms.

It works by imposing a soft-AND on the
query, where the threshold parameter is
tight (like, $q-1$ terms out of q terms of
the query should be in the doc).

NOTE

EASY TO IMPLEMENT INVERTED INDEX

Example

$Q = \text{ANTONY} \quad \text{BRUTUS} \quad \text{CAESAR} \quad \text{CALPURNIA}$

Antony	→	3 4 8 16 32 64 128
Brutus	→	2 4 8 16 32 64 128
Caesar	→	1 2 3 5 8 13 21 34
Calpurnia	→	13 16 32

COMPUTE SCORE FOR DOCS 8, 16, 32 ONLY

- Only consider high idf query terms

High inverted document frequency



short posting list



the term is not widely used

The idea is that rare terms are more valid for the search result, while more common terms contribute little to the scores (... they don't alter the rank ordering much)

NOTE EASY TO IMPLEMENT, JUST CONVNE

IN THE COSINE-SCORE'S PSEUDO CODE THE

"for t in q " in something like

"for t in $\text{high-idf}(q)$ "

• Only consider the docs in the "Champion List"

1) PREPROCESS PHASE = FIND THE m

best docs for each term (those
with higher tf-idf)

2) SEARCH PHASE = Compute cosine
between the query and those docs only
then pick the top k

NOTE works well only if m > k

NOTE m can be different between the terms

How do you retrieve the top-k documents of
AND queries?

We want a champion list with docs containing all the
terms of the query

1) Take the champion list of all the
query terms

2) Intersect them

3) Get top-k docs

NOTE: The champion lists on the single terms should contain
m > k docs or you won't return enough results.

• Fancy-hits heuristic (Champion List evolution)

1) PREPROCESS PHASE

- Assign docID by decreasing weights, using PageRank
- Sort docID by decreasing PR weights
- LITE CHAMPION LIST** • Define $FH(t) = m$ best docs for term t , using tf-idf weights
- Define $IL(t) =$ remaining docs containing the term t

QSS
The high Scorer are in FH and in the front of IL

2) SEARCH PHASE (for a Term):

- Compute the score ($+tf-idf + PR$ usually) of all the docs in $FH(t)$, like in champion lists, and keep the top K
- Scan $IL(t)$, compute the scores or above and possibly returns some docs alongside the top K .

Stopping criterion = stop after M returns

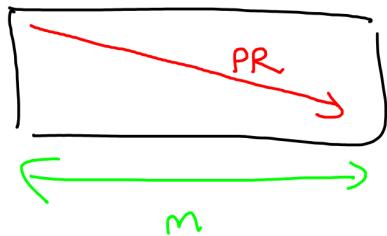
or if, given doc d :

$$tf-idf_d + PR_d < tf-idf_x \quad (? ?)$$

Example

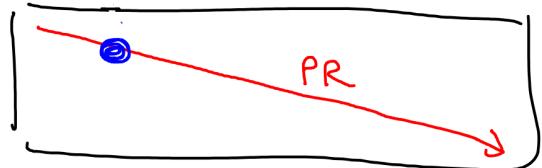
Let's say that $Q = "PISA TORRE"$,
 Let's compute FH and IL for both Term
 (1.)

PISA



Top-m by TF-IDF

FH(ϵ)



IL(ϵ)

How do you retrieve the top-K documents of AND queries?

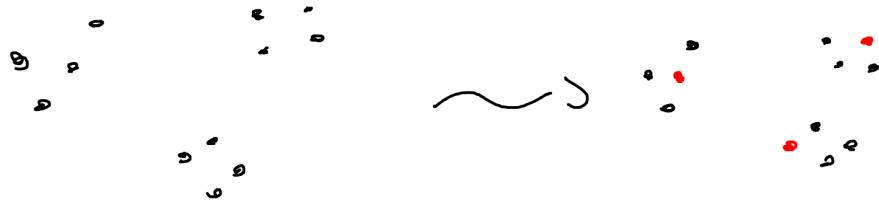
1) Take FH (include IL too?) of each term

2) Scan (in parallel?) the FHs and returns the docs where all the terms appear (kind of merge() in merge sort)

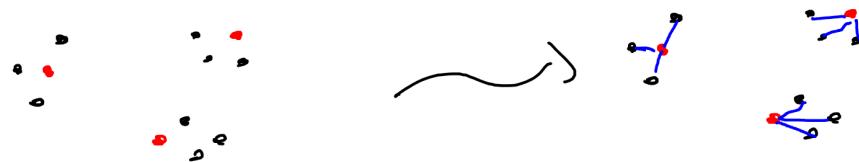
* Clustering

1) PREPROCESS PHASE :

- Pick \sqrt{N} docs at random and call them leaders



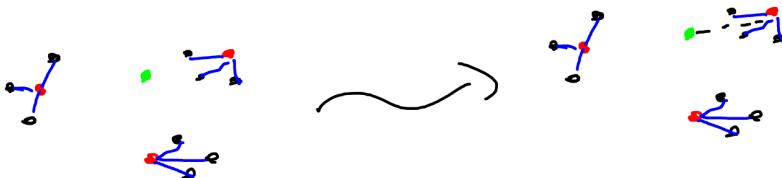
- For every other doc, locate its nearest leader and call it a follower
(i.e., each leader has $\sim \sqrt{N}$ followers)



2) QUERY PHASE :

- Given a query Q , find its nearest leader L

You could also try with more than one leader for better top-K



- Seek for the K nearest docs among L 's followers



Exact Retrieval

Given a query Q , find the exact top K docs for Q using some rating function r .
no approximations.

Simpliest strategy

- Find all documents interested by the query
- Compute score $r(d)$ for all those docs
- Sort results by score
- Return top K results

Of course, it is too costly like this and we don't have time to score every document:

we want to avoid scoring docs that won't make it into the Top K

LET'S SEE SOME
TECHNIQUES

WAND technique

A pruning method that uses a max heap over the real document scores, which is proved to contain the exact top-K docs at the end of the process.

The idea of the algorithm is similar to the Branch & Bound technique!

- Maintain a "running" threshold score Θ (which means that it changes overtime), equal to the K^{th} highest score computed so far (The "worst" top-K doc)
- Prune those docs whose scores are guaranteed to be below the threshold
- Compute exact scores for unpruned docs only

Before showing the alg...

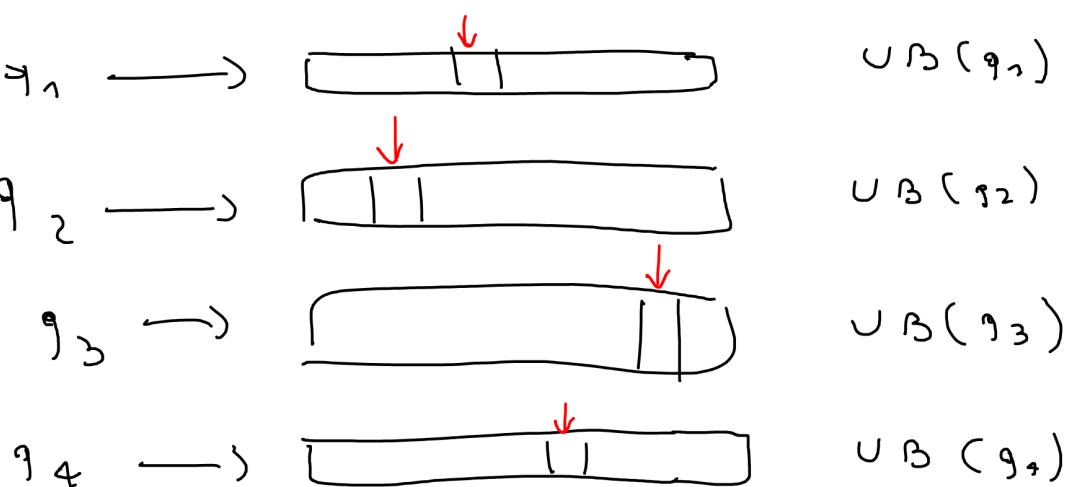
- To use WAND, we need a special iterator on the posting lists, that can only move towards the first docID larger than a certain value (since we assume ordered postings, the iterator moves only to the left, therefore towards the larger docIDs)

→ we can use skip pointers and Elias-Fano compression

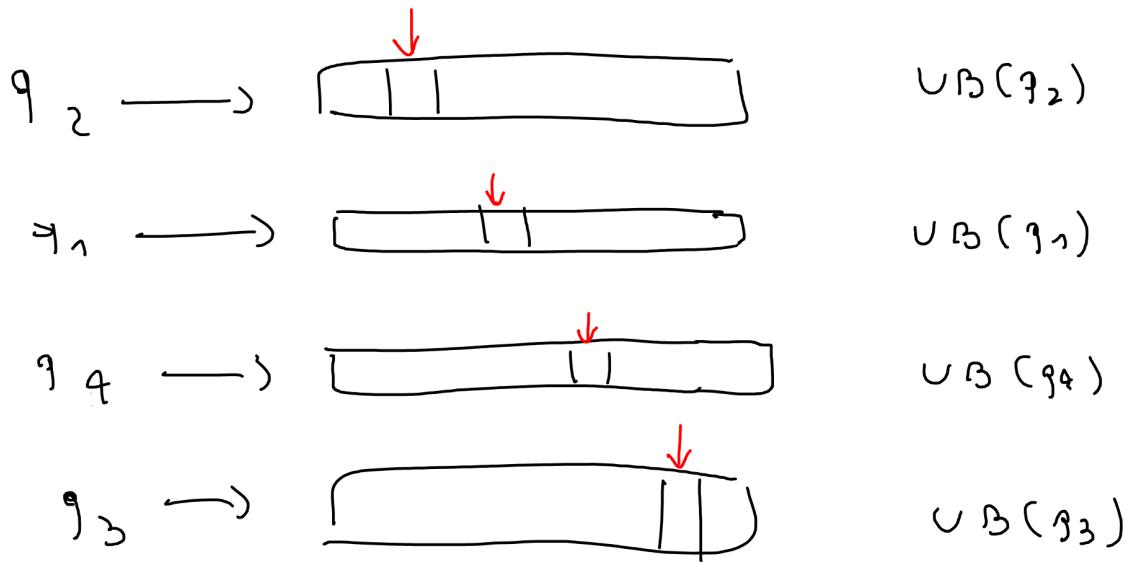
- We assume that $s_r(t, d) = \text{score of } t \text{ in } d$
- and that $s_r(d) = \text{score of } d = \sum_{t \in d} s_r(t, d)$
- every query term t have a precomputed upper bound value $UB(t)$ such that $\forall d. r(t, d) \leq UB(t)$, therefore it is the maximum possible score for t
- The threshold θ is initialized to 0, but it raises whenever the "worst" of the current top-K has a score above the threshold θ this way
- $r(d) \geq \theta$ for every d in Top-K

The alg. now

Given the query Q and the postings of the query term, let's consider a generic step of WAND, where the special position of every posting is in the following positions -



1) Sort the pointers to the postings by increasing docID



2) Find the "pivot": The first pointer in this order for which the sum of the UBs is at least equal to the threshold

$i = 0$
 $sum = -1$

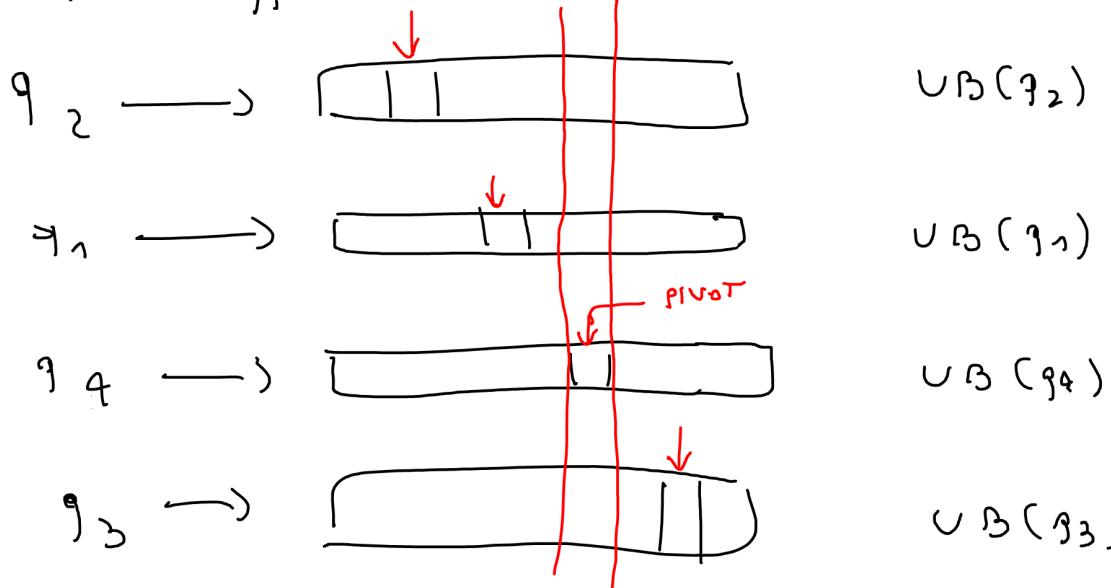
$\{q_2, q_1, q_4, q_3\}$ // sorted pointers

while $sum < \oplus$:

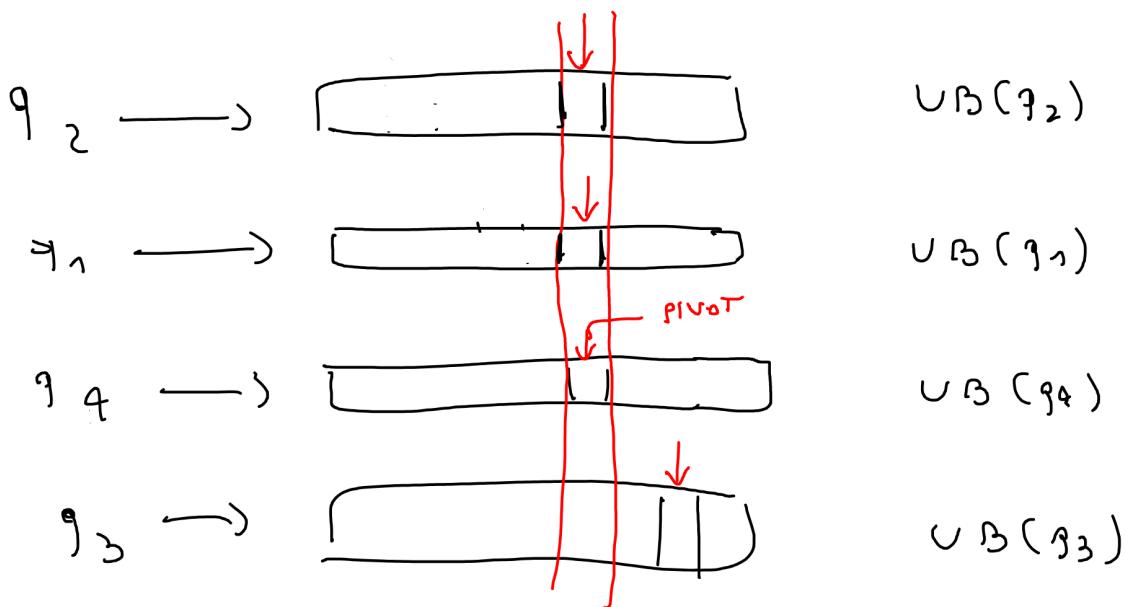
$sum += UB(q_i)$

$i++$

return q_i



3) Move special pointers to the pivot
 (this way you prune everything between the current position and the pivot)



q) If the value in the pivot is present in enough partitions (soft AND), compute the full score of that doc, else advance the pointers.

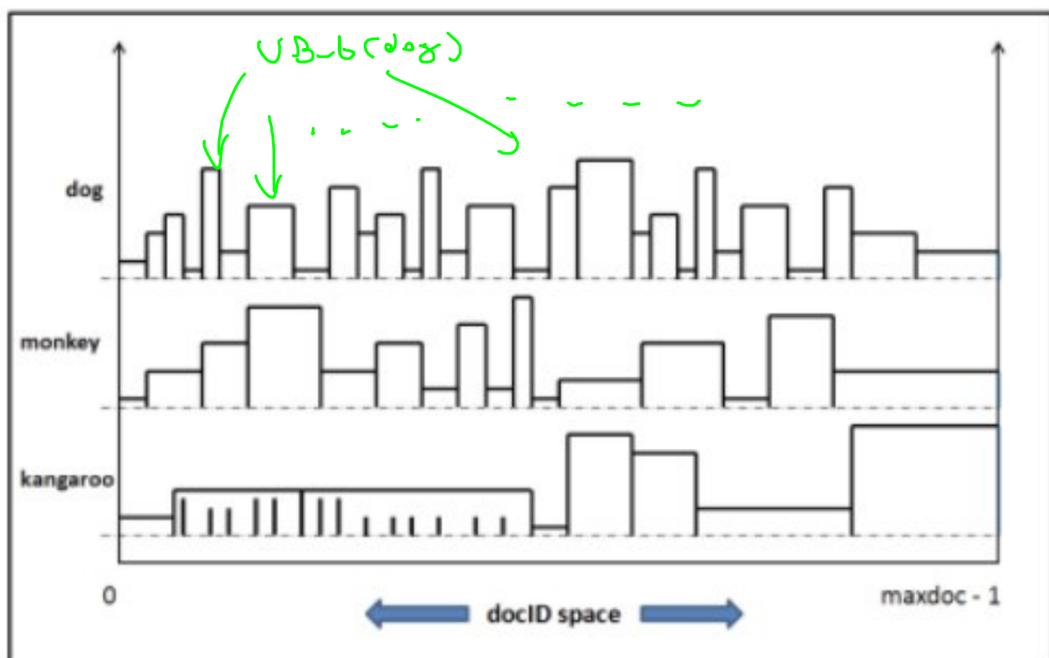
If the pivot is inserted in the current top-k ($n(d) > \theta$), update threshold θ

5) Repeat from 2)

WAND LEADS TO A 50% + REDUCTION IN SCORE COMPUTATIONS AND GIVE US A SAFE RANKING. IT PERFORMS EVEN BETTER WITH LOWER QUERIES

Blocked WAND

Like WAND, but instead of having a single UB per τ (used on the whole posting), we partition the posting into blocks and store for each block $UB_{-b(\tau)}$.



Block-Max WAND

The new version, divided in two steps:

- 1) Like WAND ($p = \text{pivot position}$
 $d = \text{pivot docID} = \text{plist}(p)$)
- 2)

Exercise

Consider the WAND alg. over 4 postings, assuming that the alg. is at the following step (iters
at the head)

$T_1 \rightarrow \dots 5, 6, 7, 8, 11$	$UB(T_1) = 0.4$
$T_2 \rightarrow \dots 2, 3, 5, 7, 8, 11$	$UB(T_2) = 2$
$T_3 \rightarrow \dots 8, 13, 15$	$UB(T_3) = 4$
$T_4 \rightarrow \dots 4, 5, 8, 3$	$UB(T_4) = 0.1$

AND $\Theta = 2, 3$

which is the next doc to be full scored?

Let's complete the next step

Step 1

1) Sort by head

$T_2 \rightarrow \dots 2, 3, 5, 7, 8, 11$	$UB(T_2) = 2$
$T_4 \rightarrow \dots 4, 5, 8, 3$	$UB(T_4) = 0.1$
$T_1 \rightarrow \dots 5, 6, 7, 8, 11$	$UB(T_1) = 0.4$
$T_3 \rightarrow \dots 8, 13, 15$	$UB(T_3) = 4$

2) find "pivot" docID

$$\text{sum} = 0 \quad \Theta = 2, 3$$

$$(\text{r2}) \quad \text{sum} = 0 + 2 = 2 \quad \Theta = 2, 3$$

$$(\text{r4}) \quad \text{sum} = 2 + 0.1 = 2.1 \quad \Theta = 2, 3$$

$$(\text{r7}) \quad \text{sum} = 2.1 + 0.9 = 2.5 \quad \Theta = 2, 3 \quad \text{so } p$$

docID = 5 passed by r4 is pivot

3) Move every one to pivot docID

$$T2 \rightarrow \dots \cancel{2, 3}, 5, 7, 8, 11$$

$$T4 \rightarrow \dots \cancel{4}, 5, 8, 3$$

$$T7 \rightarrow \dots 5, 6, 7, 8, 11$$

$$T3 \rightarrow \dots \cancel{8}, 13, 15$$

4) Does the pivot docID contains most of the terms?

$$\text{soft-AND} (\{T1, T2, T3, T4\}, 3) = \dots = 3 \quad \text{YES!!}$$

Compute the full score for docID = 5.

: if $r(5) > \Theta$, update Θ and

append docID = 5 to the Top-K results

But what if docID = 5 did not contain term T2?

soft-AND ($\{T_1, T_2, T_3, T_4\}, 3$) = ... = 2 Not enough!

Assume the pointers pointing to docID = 5
and start again

$T_2 \rightarrow \dots \cancel{2, 3, 5, 7, 8, 11}$

$T_4 \rightarrow \dots \cancel{4, 5, 8, 3}$

$T_1 \rightarrow \dots \cancel{5, 6, 7, 8, 11}$

$T_3 \rightarrow \dots 8, 13, 15$

Relevance Feedback retrieval

"Classical" Relevance feedback

The user is involved in the retrieval process of the initial set of results by giving a relevance feedback of the initial set of results by giving a relevance feedback

- 1) The user issues a start and simple query
- 2) The user marks some results as relevant or non-relevant
- 3) The system computes a better representation of the information need through the user feedback.

This way the system is able to create new query that moves towards the relevant docs and away from the irrelevant

Rocchio

$$\vec{q}_{\text{new}} = \alpha \cdot \vec{q}_{\text{old}} + \beta \frac{1}{|D_n|} \sum_{d \in D_n} \vec{d} - \gamma \sum_{d \in D_n} \vec{d}$$

\vec{q}_{old} is updated by summing the relevant documents D_r and subtracting the irrelevant document D_n , where α, β, γ are control parameters used to quantify the "trust factors"

- t) Ask the user to try the new query, and repeat
DOCUMENT RE-RANKING !!

problems

- Users often are reluctant to provide explicit feedback
- After applying the feedback (through Rochastic), it is often harder to understand why a particular document is now retrieved
- No evidence that we are not simply wasting the user's time

Pseudo - Reverse Feedback

The system assumes that the top K documents of a user's query are relevant (automating the user's job), computes reverse feedback using Rochastic and repeat.

Works well on avg, but several iterations can cause query drift

Query expansion vs Relevance feedback

In query expansion, the user gives an additional input on words and phrases, saying if a term is good or bad.

In Relevance Feedback, the user gives an additional input on documents, saying which one are relevant and which one are not.