

Assignment 3



Assignment 3

□ Documents All Pairs Similarity

↳ find all doc. that are similar at least a threshold (es. 0.95)

(useful for ex. in duplicate search
or spam removal)

□ To Be Delivered:

- Sequential implementation (Python-Numpy/Java/Scala)

- Parallel implementation

 - MapReduce / Apache Spark (Python/Java/Scala)

- Report discussing performance figures of the proposed parallel implementation

 - varying datasets (and samples), similarity thresholds
es. 0.5 or 0.95

 - varying number of workers ^{→ CPU, thread, machine}

 - max 2 pages

(LOOK AT THE PREVIOUS REPORT
COMMENTS)

Documents All Pairs Similarity

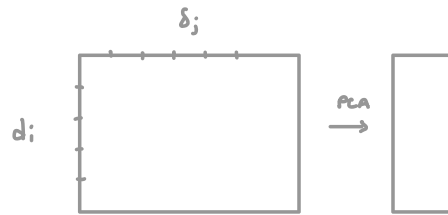
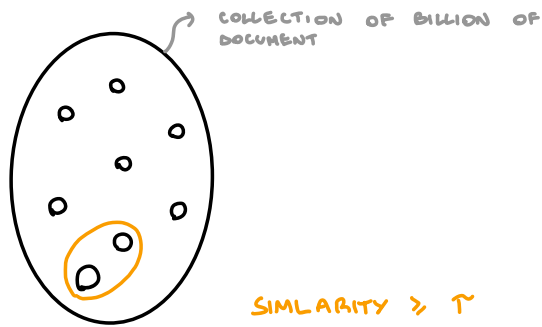
- A document is a vector d of N elements
 - N is the number of distinct words in the corpus (the lexicon)
 - $d[i]$: stores the frequency of the term i in document d ($tf(i)$).
Then d is normalized (divided by its L_2 norm)
 - additionally you can use $tf-idf$:

$$tf-idf(i) = tf(i) \cdot \ln \frac{N_{docs}}{df(i)}$$

- There are many ways to measure similarity
 - Cosine:

$$s(a, b) = \sum_{i=1 \dots N} a[i] \cdot b[i]$$

DOCUMENT ALL PAIRS SIMILARITY SEARCH



What is a document?

Can be any object \rightarrow I need to find a repr.

For textual doc I can : d_i if term not present
 $\text{term}_j \rightarrow \text{tf-idf} \sim \text{usually } \text{tf}(t) \cdot \log \left(\frac{N_f}{|d \text{ containing } t|} \right)$

For movies dataset I have : m_i u_j

Similarity can be computed using COSINE SIMILARITY that is basically a dot product

$(d_i \cdot d_j)$

if $\geq T$ doc it's append to result

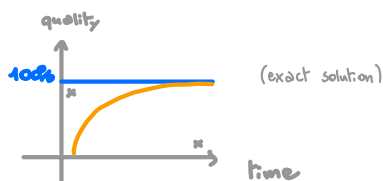
$\frac{\text{tf-idf}}{\|d_i\|_2} \sim$ to compute the score of a doc

\rightarrow

I want the EXACT similarity (but I can compare it with approximations)

- cluster docs and inside cluster do search (run alg. on clusters) (can be done on clusters in parallel)
- LSH
- PCA or JOHNSON - LINDENSTRAU

I can plot the threshold comparison



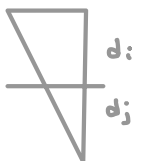
using for ex. JACCARD sim. BETWEEN RESULTS (exact and approx.)

\leftarrow (investigate the tradeoff error - speed)

COMPUTING THE EXACT SOLUTION :

For each d_i :
 For each d_j :
 compute $\text{sim}(d_i, d_j)$

$\Rightarrow O(n^2)$



THERE ARE SOME TECHNIQUES THAT SPEED UP THE COMPUTATION BUT NOT REDUCE COMPLEXITY

WE CAN ALSO SORT BY DOC LENGTH AND SKIP ALL THE DOC SHORTER THAN A THRESHOLD BECAUSE THE DOT PROD SURELY DON'T REACH THE VALUE (SKIP A LOT OF DOC IN ONE PASS WITH SORTING, RATHER THEN SKIPPING THEM ONE AT A TIME)

Sequential Algorithm

- Given the minimum required similarity *threshold*

- SIM_DOCS = 0
- For-each document d_1 in the corpus D :
 - For-each document d_2 in the corpus D :
 - if $d_1 \neq d_2$ and $s(d_1, d_2) \geq \text{threshold}$:
 - SIM_DOCS += 1

- Note: usually you are interested to the similar document pairs, rather than to the number of similar document

- Try your optimizations!!

- Datasets:

- <https://github.com/beir-cellar/beir>
- <https://grouplens.org/datasets/movielens/>

Deadline and Evaluation

- ❑ Delivery before *May 26*:
 - ❑ *+1* point if positively evaluated, *+0.5* if sufficient, re-submit if insufficient
- ❑ Or, delivery at written exam
 - ❑ +1 point if positively evaluated, +0 if sufficient, exam not passed if insufficient
- ❑ Positive Evaluation means:
 - ❑ good report, good code, good analysis