

Efficient Cosine Ranking

- Find the k docs in the corpus “nearest” to the query
 - the k largest query-doc cosines
- Efficient ranking means:
 - Computing a single cosine efficiently
 - Computing the k largest cosine values efficiently
 - Can we do this without computing all n cosines?
 - n = number of documents in corpus



Efficient Cosine Ranking

- Computing a single cosine
 - Use inverted index
 - At query time use an array of accumulators A_j to accumulate component-wise sum
 - Accumulate scores as postings lists are being processed (numerator of similarity score)

$$A_j = \sum_t (w_{q,t} w_{d,t})$$



Efficient Cosine Ranking

- For the web
 - an array of accumulators in memory is infeasible
 - so only create accumulators for docs that occur in postings list
 - dynamically create accumulators
 - put the `tf_d` scores in the postings lists themselves
 - limit docs to non-zero cosines on rare words
 - or non-zero cosines on all words
 - reduces number of accumulators



Efficient Cosine Ranking

COSINESCORE(q)

```
1  INITIALIZE( $Scores[d \in D]$ )
2  INITIALIZE( $Magnitude[d \in D]$ )
3  for each term ( $t \in q$ )
4      do  $p \leftarrow$  FETCHPOSTINGSLIST( $t$ )
5           $df_t \leftarrow$  GETCORPUSWIDESTATS( $p$ )
6           $\alpha_{t,q} \leftarrow$  WEIGHTINQUERY( $t, q, df_t$ )
7          for each  $\{d, tf_{t,d}\} \in p$ 
8              do  $Scores[d] += \alpha_{t,q} \cdot$  WEIGHTINDOCUMENT( $t, q, df_t$ )
9  for  $d \in Scores$ 
10     do NORMALIZE( $Scores[d], Magnitude[d]$ )
11  return top  $K \in Scores$ 
```



Use heap for selecting the top K Scores

- Binary tree in which each node's value $>$ the values of children
- Takes $2N$ operations to construct
 - then each of k "winners" read off in $2\log n$ steps
 - For $n = 1M$, $k = 100$ this is about 10% of the cost of sorting

