

 \mathbb{X}

 \bigcirc

Elasticsearch: Controlling relevance

 (\bigcirc)

Vitalii Melnychuk, Software Engineer



 \odot







elasticsearch



Index,

Analize





- Cluster
- Node
- Shards
- Replicas

- Index -> Table
- Document -> row
- Mapping -> schema
- Field -> column



- Document indexing
- Boolean model to find

documents

Lucene's Practical

scoring

What's next



Index time Vs Query time



• Index (noun)

Index meanings







- Index (noun)
- Index (verb)

Index meanings



- Index (noun)
- Index (verb)
- Inverted index

Index meanings



- Index (noun)
- Index (verb)
- Inverted index

Term		Doc_1		Doc_2	
Quick	1		1	х	
The	1	Х	1		
brown	1	Х	1	Х	
dog	1	Х	1		
dogs	1		1	Х	
fox	1	Х	1		
foxes	1		1	Х	
in	1		1	Х	
jumped	1	Х	1		
lazy	1	Х	1	Х	
leap	1		1	Х	
over	1	Х	1	Х	
quick	1	Х	1		
summer	1		1	Х	
the	1	Х	1		
ceap over quick summer the		x x x		× ×	



- Index (noun)
- Index (verb)
- Inverted index

Term			Doc_1		Doc_2	
(Quick	1		1	x	
	The	1	Х	Ì.		
	brown	1	х	1	х	
	dog	1	Х	1		
	dogs	1		1	Х	
	fox	1	Х	1		
	foxes	1		1	Х	
	in	1		1	Х	
	jumped	1	Х	1		
(lazy	1	Х	1	Х	
	leap	1		1	Х	
	over	1	Х	1	Х	
	quick	1	Х	1		
	summer	1		1	Х	
	the	1	Х	1		



- Document indexing
- Boolean model to

find documents

 Lucene's Practical scoring

Moving on



AND Using AND, this search would only retrieve results with Peanut Butter and Jelly.

OR Using OR, this search would retrieve results with peanut butter, with jelly, and with both. **NOT** Using NOT, this search would retrieve results with peanut butter, and exclude those with jelly or PB with jelly.



- Document indexing
- Boolean model to find

documents

 Lucene's Practical scoring

Moving on



Relevance

• Term frequency



tf(t in d) = $\sqrt{\text{frequency}}$



Relevance

- Term frequency
- Inverse document

frequency

idf(t) = 1 + log (numDocs / (docFreq + 1))







Relevance

- Term frequency
- Inverse document frequency



• Field-length norm





score (q, d) = queryNorm(q)

queryNorm = 1 / √sumOfSquaredWeights







score (q, d) = queryNorm(q) **coord(q,d)**

- Document with fox \rightarrow score: **1.5 * 1 / 3 = 0.5**
- Document with quick fox → score: **3.0 * 2 / 3 = 2.0**
- Document with quick brown fox → **score: 4.5 * 3 / 3**
 - = 4.5



```
score (q, d) =
   queryNorm(q)
   coord(q,d)
   Σ (
       tf(t in d)
       idf(t)<sup>2</sup>
       norm(t, d)
   ) (t in q)
```

tf(t in d) = $\sqrt{\text{frequency}}$

idf(t) = 1 + log (numDocs / (docFreq + 1))

norm(d) = 1 / √numTerms





```
score (q, d) =
      queryNorm(q)
      coord(q,d)
      Σ(
         tf(t in d)
         idf(t)<sup>2</sup>
         t.getBoost()
         norm(t, d)
      ) (t in q)
```

```
GET /docs 2014 */ search
"indices boost": {
 "docs 2014 10": 3,
 "docs 2014 09": 2
"query": {
 "match": {
  "title": {
    "query": "quick brown fox",
      "boost": 2
```



BM25 Algorithm

∑ (tf(t in d) * idf(t) * (k +1) /

tf(t in d) + k * (1 - b + b * termsCount/ avgTermCount) **t.getBoost()** ⊗ (t in q)



BM25's take on IDF





BM25's take on TF





TF/IDF and BM25



Frequency



BM25 Algorithm

∑ (tf(t in d) * idf(t) * (**k** +1) /

tf(t in d) + **k** * (1 - **b** + **b** * termsCount/ avgTermCount)) (t in q)



b=0.75 k=1.2

work pretty well





Scoring with scripts





```
price = doc["price"].value;
```

```
discount = doc["margin"].value;
```

score = doc[" score"];

return (price * (1 - discount)) + 1 + score;



Relevance tuning is the last 10%





Thank you!



