

Full-text search with distributed search engines

Alexander Reelsen @spinscale alex@elastic.co



Agenda

- Overview
- Indexing: Analysis, Tokenization, Filtering, on disk data structures
- Searching: Scoring, Algorithms & Optimization
- Aggregations
- Distributed systems and search
- Q & A



Filtering, on disk data structures & Optimization





Overview

Full text search introduction



Why is search so important?

Google

- Q why is search so
- Q why is search so slow in windows 10
- Q why is search so important
- Q why is search so terrible
- Q why is windows search so bad
- Q why is amazon search so bad
- Q why is outlook search so bad
- Q why is gmail search so bad
- Q why is windows search so slow
- Q why is google search so slow
- Q why is netflix search so bad

Google Search







 \mathbf{O}

Repositories	25K	25,341 repository
Code	?	
Commits	467K	elastic/elasticsearc Open Source, Distribu
Issues	180K	java search-engine
Discussions Beta	35	🖓 50.2k 👅 Java Up
Packages	30	🗐 gigablast/open-sou
Marketplace	1	Nov 20 2017 A dist
Topics	72	thter/AND. Fi ☆ 1.2k ● C++ Upda
Wikis	27K	
Users	860	A helpful checklist

Languages	
Python	4,861
Java	4,384
JavaScript	3,825
HTML	1,919
PHP	1,474
C++	830
CSS	754
C#	555
Jupyter Notebook	483
Ruby	477

Advanced	search	Cheat	sheet
avanoca	0001011	onour	011000

	Nov 20 Intel/AN	2017 MD. Fr	A distrik
	☆ 1.2k	• C++	Update
Ļ	marco QA he	biederm Ipful che	o <mark>ann/se</mark> cklist / c
	seo	keyword	optim
	公 1.5k	● HTM	L Upda
Ļ	subins	2000/s	earch
	An Ope	n Source	e Search
	php	search	search
	☆ 119	PHP	Update
Ļ	mtiany	/an/Funp	oySpide

借鉴自慕	课网	[Scra	apy 1
elastics	earch-	analy	sis-ik
zhihu	scrap	у	lagou
分726	🔵 Pyt	hon	Apad

Ļ	NeXTs/Jets.js Native CSS search eng		
	css search-engine		
	☆ 2.8k	JavaScript	



e \sim Marketplace Pricing \sim	search engine	Sign in Sig	gn up
25,341 repository results	Sort: Best match	-	
 □ elastic/elasticsearch Open Source, Distributed, RESTful Search Engine java search-engine elasticsearch ☆ 50.2k ● Java Updated 3 minutes ago 213 issues need help 			
 ☐ gigablast/open-source-search-engine Nov 20 2017 A distributed open source search engine and spide Intel/AMD. Fr ☆ 1.2k ● C++ Updated on 4 May 	er/crawler written in C/C++ for Linux o	n	
 □ marcobiedermann/search-engine-optimization ○ A helpful checklist / collection of Search Engine Optimization (See keyword optimization search google awesome conditions) ○ keyword optimization search google awesome condition ○ 1.5k ● HTML Updated 7 days ago 	SEO) tips and techniques. Illection engine tips awesome-list		
 □ subins2000/search An Open Source Search Engine php search search-engine crawler ☆ 119 ● PHP Updated on 24 Jul 2017 			
 □ mtianyan/FunpySpiderSearchEngine 借鉴自慕课网【Scrapy 1.6.0爬取数据 + ElasticSearch6.8.0+Django2 elasticsearch-analysis-ik mysql redis django elasticsearch zhihu scrapy lagou ☆ 726 ● Python Apache-2.0 license Updated on 5 Apr 	2.2搜索引擎】 python search-engine spider		
 ReXTs/Jets.js Native CSS search engine css search-engine ☆ 2.8k ● JavaScript MIT license Updated on 7 Nov 2019 			





elast



Sign up

Hot Network Questions



WOMEN | MEN | KIDS

Get the Look NEW Clothing Shoes

Q hemd grün



Clothing

Shoes

Sport

Accessories

Premium

Care

Gift vouchers

Sale

'hemd grün'× Sort ∽ Size ∽ Brand ∽ Colour



GANT GANT HERREN HEMD REGULAR ... 99,99€

elastic

33,334

zalando	English	C	Wish list	H Your bag	
Price ~ Material ~ Delivery ~ 2 items					
<image/>					
-56%					
Andrew James 35,99 € ANDREW JAMES HEMD - Shirt - gr 15,99 €					
< Page 1 of 1 >					



40 Eigentumswohnungen in Schwabing Suche speichern





















	← All Shared C	Content Sources		
	<table-cell-rows> Jira</table-cell-rows>	1		Source Overview
643	Connector	Created		Content Summary
4	Jira	July 29, 2019		CONTENT TYPE
සී	Overview		>	Story
R	Content			Project
Å				Other
•••	Remove Ji	ira		Total Documents
$\langle \mathcal{O} \rangle$				Recent Activity
N				EVENT
E)				Syncing
				Sync
				Sync
				Created



,	
	Manage
	ITEMS
	42
	4
	89
	135

Recent Activity

TIME

Product

Design

Less than a minute ago

1 day ago

1 day ago

1 day ago

















SELECT * FROM products WHERE name LIKE = '%topf%'





grep "topf" my_dataset.txt





Problem

- Scales linearly with the data set size
- Relevancy
- Spell correction
- Synonyms
- ases









Inverted Index





The quick brown fox jumped over the lazy dog





The quick brown fox jumped over the laz dog









The brown dog fox jumped lazy over quick the







Quick brown foxes leap over lazy dogs



in summer





lazy dog



lazy AND dog

[1,2] AND [1] = [1]



lazy OR dog

[1,2] OR [1] = [1,2]



Technologies used today

- Apache Lucene (search library)
- Elasticsearch (distributed search engine built on top of Apache Lucene)









Indexing

Analysis, Tokenization, Filtering Data structures



quick



elastic

quick



Analysis: Tokenizer & Token Filters











quick brown fox







quick brown fox







quick brown fox the lazy, white dog.





quick brown fox the lazy white dog





quick_brown_fox the_lazy_white_dog

Unicode® Standard Annex #29

UNICODE TEXT SEGMENTATION

Summary

This annex describes guidelines for determining default segmentation boundaries between certain significant text elements: grapheme clusters ("user-perceived characters"), words, and sentences. For line boundaries, see [UAX14].

https://unicode.org/reports/tr29/





quick_brown_fox
the_lazy_white_dog
https://www.jade-hs.de




Tokenization

quick_brown_fox
the_lazy_white_dog
https_www.jade_hs.de











Quick The brown dog dogs fox foxes in jumped lazy leap over quick summer the elastic



The Quick brown fox





Lowercase

The Quick brown fox the quick brown fox





Lowercase Stopwords



quick brown fox



Lowercase S

The Quick brown fox the quick brown fox



Stopwords



quick orown Fox

quick, fast brown fox



Stopwords Lowercase Synonyms

Tokens can be changed, added, removed



quick, fast brown fox



Lowercase S



Queries need to be processed as well!



Stopwords Synonyms

quick orown fox

quick, fast brown fox



More analysis strategies

- Phonetic analysis: Meyer vs. Meier
- Stemming: foxes -> fox
- Compounding: Blumentopf → blumen topf











(On-Disk) Data structures



What else is in an inverted index?

- Documents: Find documents
- Term frequencies: Relevancy
- Positions: Positional Queries
- Offsets: Highlighting
- Stored fields: The original data





Segment: Unit of work

- A fully self sufficient inverted index
- An index consists of a number of segments
- New segments are created for newly added documents
- Segments are immutable!









Read-only data structures

- Pro: Write-once, sequentially
- Pro: Lock-free reading
- Pro: File system cache
- Contra: in-place updates & deletes
- Contra: Housekeeping
- Contra: Transactions elastic



Segment: Deletes

- Mark a document as deleted in a special file
- Exclude it from searches
- No space is freed!





Segment: Merging

- Number of segments needs to be kept reasonable
- Merge multiple segments into one (smaller index)
- Delete expired documents







Segment: Merging

- Number of segments needs to be kept reasonable
- Merge multiple segments into one (smaller index)
- Delete expired documents















Searching

Precision vs. recall Scoring Algorithms and optimizations









Relevancy

- Textbook answer: How well matches a document a query?
- Business answer: Are the top search results those that make me the most money?
 - marketplace
 - hotel booking website
 - newspaper website



matches a document a query? It search results those that









Scoring: lazy dog

- Naive: increase a counter if a term is matched
- "the lazy dog" => score 2
- "the lazy frog" => score 1
- "the lazy lazy lazy lazy cat" => score 4 or 1?







Scoring: More than term frequency

- How about incorporating information about the whole document corpus in scoring?
- Are lesser common terms more relevant?
 news paper: "dieselgate news"





Scoring: TF-IDF

- a field
- Inverse document frequency: inverse function of the number of documents in which it occurs



Term frequency: number of times a term occurs in



Scoring: Vector space model

- Each term is a dimension
- The length is based on tf-idf calculation
- Similarity is the angle between vectors
- Cosine similarity: best match == angle 0°





Scoring: TF-IDF in Lucene

score(q,d) = $\sum (tf(t in d) \cdot idf(t)^2 \cdot t.getBoost() \cdot norm(t,d))$







Lucene's Practical Scoring Function is derived from the above. The color codes demonstrate how it relates to those of the conceptual formula:

score(q,d) =
$$\sum (tf(t \text{ in } d) \cdot idf(t)^2 \cdot t.get)$$

t in q

Lucene Practical Scoring Function

where

1. *tf(t in d)* correlates to the term's *frequency*, defined as the number of times term *t* appears in the currently scored document *d*. Documents that have more occurrences of a given term receive a higher score. Note that *tf(t in q)* is assumed to be *1* and therefore it does not appear in this equation, However if a query contains twice the same term, there will be two term-queries with that same term and hence the computation would still be correct (although not very efficient). The default computation for *tf(t in d)* in ClassicSimilarity is:

```
tf(t in d) = frequency^{\frac{1}{2}}
```

2. *idf(t)* stands for Inverse Document Frequency. This value correlates to the inverse of *docFreq* (the number of documents in which the term t appears). This means rarer terms give higher contribution to the total score. *idf(t)* appears for t in both the query and the document, hence it is squared in the equation. The default computation for *idf(t)* in ClassicSimilarity is:

$$doct$$

 $idf(t) = 1 + log(-----doct$

- 3. **t.getBoost()** is a search time boost of term t in the query q as specified in the query text (see query syntax), or as set by wrapping with BoostQuery. Notice that there is really no direct API for accessing a boost of one term in a multi term query, but rather multi terms are represented in a query as multi TermQuery objects, and so the boost of a term in the query is accessible by calling the sub-query getBoost().
- 4. norm(t,d) is an index-time boost factor that solely depends on the number of tokens of this field in the document, so that shorter fields contribute more to the score.



```
Count+1
_____
```

cFreq+1





- Default in Apache Lucene/Elasticsearch
- Works better with stopwords (high TF)
- Term frequency saturation
- Improved field length normalization (per document)









https://www.elastic.co/guide/en/elasticsearch/guide/2.x/pluggable-similarites.html



requency







Precision vs. recall



Precision and Recall





Precision and Recall









Precision and Recall







True positives







True negatives

relevant documents







False positives

relevant documents






False negatives

relevant documents







Precision and recall

 Precision: How many selected documents are relevant?

 Recall: How many relevant documents are selected









Under the hood



Optimizations everywhere

- leap frogging, skip lists
- top-k
- two phase iterations
- integer compression











Query two phase iteration



Two phase iteration: Phrase query

- Phrase query: "quick fox"
- Approximation phase: document contains terms quick and fox
- Verification phase: read positions of terms





Two phase iteration: Geo distance query

- Geo distance query: Distance from reference point
- Approximation phase: bbox around point
- Verification phase: exact distance calculation





Two phase iteration: Geo distance query

```
GET /my_locations/_search
"query":
  "bool" :
    "filter" :
          "lon" : -70
```



"geo_distance" : {
 "distance" : "200km",
 "pin.location" : {
 "lat" : 40,



Two phase iteration: several queries

- Powerful when several queries are used
- "quick fox" AND brown
- Approximation: quick AND fox AND brown
- Verification: "quick fox" position check for hits





Skip lists & leap frogging





Skip lists

- Term dictionary is a sorted skip list
- Skip list is a linked list with 'express lanes' to leap forward





https://en.wikipedia.org/wiki/Skip_list



Leap frogging



elasticsearch AND kibana AND logstash



Leap frogging elasticsearch AND



elasticsearch AND kibana AND logstash









Ę-





Ę-







Ę-







Ę-



266
102
98
60
18
5
1







Ę.



	266
	102
	98
	60
	18
)	5
	1





Ę-



266
102
98
60
18
5
1







266
102
98
60
18
5
1







266
102
98
60
18
5
1







266
102
98
60
18
5
1







Ser.

	266	
Ī	102	Hit!
	98	
	60	
	18	
	5	Hit!
	1	









Ser.

	266	
Ī	102	Hit!
	98	
	60	
	18	
	5	Hit!
	1	









Ser.

	266	
Ī	102	Hit!
	98	
	60	
	18	
	5	Hit!
	1	







Ser.

	266	
Ī	102	Hit!
	98	
	60	
	18	
	5	Hit!
	1	







Ser.

	266	
Ī	102	Hit!
	98	
	60	
	18	
	5	Hit!
	1	







Top-k retrieval



Top-k retrieva

- elasticsearch OR kibana
- top 10 results wanted
- maximum score for kibana is 3.0
- maximum score for elasticsearch is 5.0
- collecting documents: when 10th hit has score > 3, then only documents with elasticsearch need to be collected
- total hit count is not accurate





Top-k retrieval



Ę. elastic

Oct 2017

14

12

10

8

BooleanQuery (OR, high freq, medium freq term)

							2017/12	/18 09:41:42 31	:
				CJ			Q D . 5		P
				CK		who man	many CI	N COmm	çq
							CM	a differ the contract	
							W		
	CD CE	СССССН							
		and the second second	~~~~						
Jan	2018 Apr	2018 Jul 2	2018	Oct Date	2018 Ja	n 2019	Apr 2019	Jul 20)19









Order index by field values

each segment is sorted before write

criteria can be chosen by the user

5 | 2 | 3 | 1 | 4

retrieve 5 | 2 | 3 | 1 | 4





top 2 sort 5 | 4 | 3 | 2 | 1 5 | 4



Order index by field values

- each segment is sorted before write
- criteria can be chosen by the user



early termination 5 | 4









Aggregations

Reducing data



Aggregations

	\boldsymbol{r}	
		,
-	-	

nike hoodie

Inspiration	Damen .	Herren	Kinder	
Multimedia	Haushalt	Küche	Heimtext	ili

Startseite | Suchergebnis für nike hoodie (129)

^	
	000
	5/62
^	ATTE (
	Cart - V
^	
	NIKE SPORTSWEAR Nike Sportswear Kapuzensweatshirt
^	»M NSW CLUB HOODIE PO BB GX< € 54,99
-	





Sortieren nach Topseller



NIKE SPORTSWEAR Nike Sportswear Kapuzensweatshirt »M NSW CLUB HOODIE PO BB« € **49,99**





\$|

NIKE SPORTSWEAR

Nike Sportswear Kapuzensweatshirt »BOYS NIKE SPORTSWEAR HOODIE CLUB FLEECE BRUSEHD«

€ 39,99



> Mehr aus der Serie
























Sneakers



Oxfords



Sneakers

Boots



Sneakers











Pumps



Oxfords





Sneakers











 $\langle \cdot \rangle$ elastic



Pumps





Oxfords



bucket agg

metric agg

doc_count



Sneakers



















Pumps









23

bucket agg

metric agg

avg price



Aggregations

- bucket: terms, histogram, geo, range, sampler, significant text, nested
- metric: value_count, avg, min, max, sum, stats, median deviation, geo, percentile, cardinality,
- differencing



 pipeline: min, max, sum, avg, derivative, stats, percentiles, cumulative sum, moving average, moving function, serial





Distributed systems & search

Fanning out a search, reducing the results







Elasticsearch







•



Elasticsearch in 10 seconds

- Search Engine (FTS, Analytics, Geo), near real-time
- Distributed, scalable, highly available, resilient
- Interface: HTTP & JSON
- Centrepiece of the Elastic Stack
- Uneducated conservative guess: Tens of thousands of clusters worldwide, hundreds of thousands of instances









Distributed systems



Distributed systems

- How do nodes communicate with each other?
- Who is taking and executing decisions?
- Failure detection?
- Replication strategy?
- Consistency?
- Enter consensus algorithms... elastic







A fundamental problem in distributed computing and multi-agent systems is to achieve overall system reliability in the presence of a number of faulty processes. This often requires processes to agree on some data value that is needed during computation

https://en.wikipedia.org/wiki/Consensus_(computer_science)

Consensus algorithms

- Leader based: Paxos, Raft
- Non leader based: BTC, gossip





- Custom consensus algorithm, improving the existing one
- Formally verified
- Optimized for Elasticsearch use-case (rolling) restarts, growing/shrinking clusters, log-ofoperations vs. cluster state)































































Master node tasks

- Deciding where data should be stored
- Pinging other nodes
- Reacting on node leaves/joins
- Updating cluster state
- Distributing cluster state















































































Distributed search





Primary Shard (Lucene Index)











Replica shard (copy)

















node 3

p1

Replica shard (copy)

r1

node 4



- Shard: Lucene index, unit of scale
- Primary shard: Write scalability
- Replica shard: Read scalability, availability















1. Client connects any node with search request



2. Execute query against shards









nects any node with search request



3. top-k search results are returned to coordinating node



1. Client con





nects any node with search request



4. Create real top-k result list













Distributed search in Elasticsearch 5. Fetch original documents









nects any node with search request


Distributed search in Elasticsearch



















Aggregations - cardinality









POST /sales/_search?size=0 "aggs" : { "type_count" : { "cardinality" : { "field" : "type"



Aggregations - cardinality



25



40





How many distinct elements are in my index?

What is the total? 40? 65?

Naive solution: merge data to single dataset and count. Doesn't scale!

Solution: Use HyperLogLog++



HyperLogLog++

- Hash based counting
- Trades in memory for accuracy
- Fixed memory usage, based on configurable precision
- Result: Small mergeable data structure, can easily be sent over the network elastic



Aggregations - percentile







GET latency/_search "size": 0, "aggs" : { "load_time_outlier" : { "percentiles" : { "field" : "load_time"



T-Digest

- Percentiles are divided into buckets
- When buckets grow over a boundary, approximation kicks in, saving memory in the process
- The exact level of inaccuracy is difficult to generalize
- Alternative: HDR histograms



Extreme percentiles are more accurate than the Median



Probabilistic data structures

- bloom/cuckoo/quotient filters (membership check)
- HyperLogLog++ (cardinality)
- T-Digest, DDSketch, HDR histogram (percentile)
- Count-Min sketch (frequency, top-k)
- Hashing (similarity)







Demo

Try it out yourself!

https://ela.st/jade-hochschule-samples





Elastic Cloud Free 30 day trial

elastic 🔗

Less

We hope you learned something new -- put your knowledge to the test and try out the Elastic Stack today.

- ✓ 30-day free trial
- No credit card required
- Get the latest versions, powerful features, and

Enter your email

Start Free Trial

https://ela.st/hack-your-future-2020

Deploy Elasticsearch and Kibana in 3 Minutes or

optimized deployment templates for your use case.

Deployments Custom pluging Account Help

Create deployment Name your deployment Give your deployment a name 2 Select a cloud platform aws Amazon Web Services 3 Select a region US East (N. Virginia) Asia Pacific (Singapore) EU (Frankfurt) 4 Set up your deployment Elastic Stack version 6.5.1 Edit

Optimize your deployment





Upcoming trends & summary

... or why you should take a closer look at search



Search is not just google...

- "Just google it" does not cut it
- Enterprise search: Intranet/G-Drive/Dropbox
- Ecommerce search
- SIEM
- Observability: Logging, APM & Metrics





Search is not 'done'

- Constant improvement
- Data structures & algorithms (BKD tree for geo shapes)
- Academic research moves to industry thanks to Apache Lucene





Search is still tough

- Language specific analysis
- Smart query parsing (nike red hoodie x1)
- Geo based search
- Anomaly detection
- Incoporating feedback loops









Upcoming trends

- Learning-to-Rank
- Deep Learning
- Feedback loop





Summary

- Everything is a search problem!
- Search is hard... and interesting
- Distributed systems are hard... and interesting
- Domain knowledge required
- Data keeps exploding, good job chances!







Literature

Books, books, books



classification Christopher D. Manning search Prabhakar Raghavan **Hinrich Schütze** recision links spam Introduction to Information

Retrieval

query

With applications for Solr and Elasticsearch

Relevant SEARCH

John Berryma

O'REILLY°



Elasticsearch The Definitive Guide

A DISTRIBUTED REAL-TIME SEARCH AND ANALYTICS ENGINE

INFORMATION RETRIEVAL

Implementing and Evaluating Search Engines

Stefan Büttcher Charles L. A. Clarke Gordon V. Cormack Covers Apache Lucene 3.0

Lucene **SECOND EDITION**

Michael McCandless Erik Hatcher Otis Gospodnetić FOREWORD BY DOUG CUTTING

Ian H. Witten Alistair Moffat Timothy C. Bell

EEP LEARNING for Search

Managing Compressing and Indexing Documents and Images Gigabytes

SECOND EDITION



Tommaso Teofili Foreword by Chris Mattmann

MANNING





Resources

Links, links, links



Links

- https://lucene.apache.org/core/8_2_0/core/org/apache/lucene/search/similarities/ **TFIDFSimilarity.html**
- https://www.elastic.co/blog/whats-new-in-lucene-8
- https://www.elastic.co/blog/faster-retrieval-of-top-hits-in-elasticsearch-with-block-max-wand
- https://speakerdeck.com/elastic/amusing-algorithms-and-data-structures
- https://www.elastic.co/blog/index-sorting-elasticsearch-6-0
- https://raft.github.io/
- https://github.com/elastic/elasticsearch-formal-models https://gist.github.com/spinscale/b62c8b357fae7db3f14b7d3127758951





Links - probabilistic data structures

https://github.com/addthis/stream-lib https://github.com/DataDog/sketches-java https://github.com/HdrHistogram/HdrHistogram https://github.com/JohnStarich/java-skip-list https://github.com/addthis/stream-lib https://static.googleusercontent.com/media/research.google.com/fr/ pubs/archive/40671.pdf







Q & A

Alexander Reelsen alex@elastic.co @spinscale

