# Elasticsearch in Ecommerce

Alexander Reelsen

@spinscale

alexander.reelsen@elasticsearch.com

elasticsearch.

# about

- ## Me

  Interested in metrics, ops and the web
  Likes the JVM
  Working with elasticsearch since 2011

- ## Elasticsearch, founded in 2012

  Products: Elasticsearch, Logstash, Kibana
  Professional services: Support & development subscriptions
  Trainings

**elasticsearch.**

# Agenda

- Introduction

- Ecommerce Use-Cases
  Product/Full-text search
  Logfiles
  Analytics

- Elasticsearch 1.0

- Q & A

elasticsearch.

# Introduction

elasticsearch.

# Unstructured search

elasticsearch.

# Structured search

elasticsearch.

# Enrichment

elasticsearch.

# Sorting

**GitHub**

Explore  Features  Enterprise  Blog

Sign up    Sign in

Search

elasticsearch                                              Search

Sort: Best match ▾

We've found 317 repository results

📖 Repositories        317
<> Code              17,981
① Issues              2,008
👤 Users                  2

**elasticsearch/elasticsearch**                    Java  ★ 4,683  ⑂ 1,097
Open Source, Distributed, RESTful Search Engine
Last updated 2 hours ago

Languages

Java                    ⊗
Ruby                   167
JavaScript             139
Python                 117
PHP                     69
Shell                   49
Puppet                  40
Perl                    38
Scala                   16
C#                      13

**richardwilly98/elasticsearch-river-mongodb**     Java  ★ 308  ⑂ 48
MongoDB River Plugin for ElasticSearch
Last updated 2 minutes ago

**jprante/elasticsearch-river-jdbc**               Java  ★ 170  ⑂ 70
JDBC river for Elasticsearch
Last updated 12 days ago

**elasticsearch/elasticsearch-hadoop**             Java  ★ 79  ⑂ 28
Read and write data to/from ElasticSearch within Hadoop
Last updated 3 days ago

**elasticsearch.**

# Pagination

elasticsearch.

# Aggregation

elasticsearch.

# Suggestions

elasticsearch.

# Elasticsearch in 10 seconds

- Schema-free, REST & JSON based distributed document store

- Open Source: Apache License 2.0

- Zero configuration

- Written in Java, extensible

elasticsearch.

# Zero configuration

```
$ wget https://download.elasticsearch.org/...

$ tar -xf elasticsearch-1.0.0.RC1.tar.gz

$ ./elasticsearch-1.0.0.RC1/bin/elasticsearch -f

...
[2014-01-19 14:53:11,508][INFO ][node] [Scanner] started
...
```

elasticsearch.

# Is it alive?

```
» curl localhost:9200
{
  "status" : 200,
  "name" : "Scanner",
  "version" : {
    "number" : "1.0.0",
    "build_hash" : "e018cda7e7a32643d59e0ac3cdb412ccc239af04",
    "build_timestamp" : "2014-01-17T15:11:47Z",
    "build_snapshot" : true,
    "lucene_version" : "4.6"
  },
  "tagline" : "You Know, for Search"
}
```

elasticsearch.

# Create...

```
» curl -XPUT localhost:9200/books/book/1 -d '
{
  "title" : "Elasticsearch - The definitive guide",
  "authors" : "Clinton Gormley",
  "started" : "2013-02-04",
  "pages" : 230
}'
```

elasticsearch.

# Update…

```
» curl -XPUT localhost:9200/books/book/1 -d '
{
  "title" : "Elasticsearch - The definitive guide",
  "authors" : [ "Clinton Gormley", "Zachary Tong" ],
  "started" : "2013-02-04",
  "pages" : 230
}'
```

elasticsearch.

# Delete…

```
» curl -X DELETE localhost:9200/books/book/1
```

# Realtime GET…

```
» curl -X GET localhost:9200/books/book/1
» curl -X GET localhost:9200/books/book/1/_source
```

elasticsearch.

# Search

```
» curl –XGET localhost:9200/books/_search?q=elasticsearch
```

```
{
  "took" : 2, "timed_out" : false,
  "_shards" : { "total" : 5, "successful" : 5, "failed" : 0 },
  "hits" : {
    "total" : 1, "max_score" : 0.076713204,
    "hits" : [ {
      "_index" : "books", "_type" : "book", "_id" : "1",
      "_score" : 0.076713204, "_source" : {
        "title" : "Elasticsearch — The definitive guide",
        "authors" : [ "Clinton Gormley", "Zachary Tong" ],
        "started" : "2013–02–04", "pages" : 230
      }
    } ]
  }
}
```

**elasticsearch.**

# Search - Query DSL

```
» curl -XGET 'localhost:9200/books/book/_search' -d '{

    "query": {
        "filtered" : {
            "query" : {
                "match": {
                    "text" :  {
                        "query" : "To Be Or Not To Be",
                        "cutoff_frequency" : 0.01
                    }
                }
            },
            "filter" : {
                "range": {
                    "price": {
                        "gte": 20.0
                        "lte": 50.0
            ...
        }
}'
```

elasticsearch.

# Distributed & scalable

- ## Replication
  Read scalability
  Removing SPOF

- ## Sharding
  Split logical data over several machines
  Write scalability
  Control data flows

elasticsearch.

# Distributed & scalable

## node 1

### orders

| | |
|---|---|
| 1 | 2 |
| 2 | 4 |

### products

| | |
|---|---|
| 1 | 2 |

```
curl -X PUT localhost:9200/orders -d '{
    "settings.index.number_of_shards" : 4
    "settings.index.number_of_replicas" : 1
}'
```

```
curl -X PUT localhost:9200/products -d '{
    "settings.index.number_of_shards" : 2
    "settings.index.number_of_replicas" : 0
}'
```

**elasticsearch.**

# Distributed and scalable

**node 1**

orders
- 1
- 2
- 3
- 4

products
- 1

**node 2**

orders
- 1
- 2
- 3
- 4

products
- 2

elasticsearch.

# Distributed & scalable

**node 1**

orders

1     2

     4

products

1

**node 2**

orders

2

3

products

2

**node 3**

orders

1

3     4

products

elasticsearch.

# Distributed & scalable

- JVM (high level & high performance if done right)

- Netty (async networking on top of the JVM)

- Lucene (fulltext search library)

- HPPC (high performance primitive collections)

- Google Guice (for extension & dependencies)

elasticsearch.

# A request under the hood



REST Event Loop

Transport Event Loop

Action Event Loop

Request

Response

elasticsearch.

# Think async!

- Enforces event driven architecture

- Support for non-blocking model

- Enforce loose coupling

- Prefers push over pull

- Callback based concurrency

- Helps to avoid contention on resources / threads

elasticsearch.

# Ecosystem

- Plugins

- Clients for many languages
  Ruby, python, php, perl, javascript, (.NET coming)
  Scala, clojure, go

- Kibana

- Logstash

- Hadoop integration

elasticsearch.

# Use-case:
# Product search engine

elasticsearch.

# Product search engine

- Just index all your products and be happy?
  Search is not that easy

- Gathered experience at an b2b ecommerce platform in the hotel and gastronomy sector
  First solution was self written using bobo/zoie turned out to be unmaintainable
  Switched to elasticsearch then

- Decompounding, Suggestions, Faceting, Custom scoring, Analytics, Price agents, Query optimization, beyond search

elasticsearch.

# Domain specific knowledge

- ## Search term: Topf
  What is expected? Blumentopf? Kochtopf?
  Or: Tuch (Handtuch, Halstuch, Geschirrtuch)
  Or: Decke (Tischdecke, Löschdecke, Mitteldecke)

- ## Decompounding (compound word token filter)
  Blumentopf also needs to match Leuchtblumentopf

- ## Synonyms
  Portmonee/Portemonnaie/Geldbörse

elasticsearch.

# Neutrality? Really?

- Is full-text search relevancy really your preferred scoring algorithm?

- Possible influential factors
  Age of the product, been ordered in last 24h
  On stock?
  Provision
  No shipping costs
  Special offer
  Rating (product or seller)

  http://www.elasticsearch.org/guide/en/elasticsearch/reference/current/query-dsl-function-score-query.html

elasticsearch.

# Faceting & Filtering

- ## Products grouped by
  Category
  Material
  Brand

- ## Allowing to filter
  All of the facets
  Price range
  Color
  Seller
  Ratings (hard!)

**Kategorien**
**Elektronik & Foto**
Fernseher
+ Mehr...
+ Alle 35 Kategorien

**Versandoption** (Was ist das?)
☐ Kostenlose Lieferung ab EUR 20 Bestellwert

**Displaygröße von Fernsehern**
☐ 51 cm (20") & kleiner
☐ 53 - 59 cm (21-23")
☐ 61 - 76 cm (24-30")
☐ 79 - 99 cm (31-39")
☐ 102 - 114 cm (40-45")
☐ 116 - 120 cm (46-47")
☐ 121 - 140 cm (48-55")
☐ 142 cm (56") & mehr

**Fernseher-Funktionalität**
☐ Smart / Internet
☐ 3D
☐ HbbTV

**Farbe**

**Fernseher-Seitenverhältnis**
☐ 16:9 Wide screen
☐ 4:3 Standard

**Displaytechnologie von Fernsehern**
☐ LED Backlight
☐ LCD
☐ Plasma

**Durchschn. Kundenrezension**
★★★★☆ & mehr
★★★☆☆ & mehr
★★☆☆☆ & mehr
★☆☆☆☆ & mehr

elasticsearch.

# Product variants?

- How to handle product variants?
  Same product by the same merchant
  Same product in different sizes, colours (clothing)

- Solution: Patched elasticsearch with grouping support, which was done by creating an image hash from the image and grouping on it

- Unsolved: Same product by different merchant
  Unless the exact same image is used, unlikely

- Better solution: Parent/child support

elasticsearch.

# Notification with Percolation

- Customer: If a product matches name *X* and costs below price *Y*, is color *Z*, then I want to get a mail
  More likely: Notify customer, when it is back on stock

- Enter percolation!
  Not: Index a document and fire a query
  But: Index a query and check a document against if it matches

  https://speakerdeck.com/javanna/whats-new-in-percolator

elasticsearch.

# More than pure search

- ## Users (ab)use the search bar for everything
  Imprint, Careers, Jobs, special offer
  Requires a special component between web app and search
  which redirects for special search terms to landing pages

- ## Analytics
  Save all your queries, and analyze
  Most searched terms
  Most searched terms with zero results
  Searched terms, which lead to an add-to-cart action
  Searched terms, which lead to complete abort

elasticsearch.

# Beware: Data quality

- Data quality can kill all your search improvements in no time

  Tough bet, if you rely on external products
  Will require you to have an own ETL pipeline, before the data goes into search or your platform (hard!)

- Less products, but more enriched results in more relevant searches

- Tough in a multi merchant environment in a non IT driven industry with lots of small businesses

elasticsearch.

# Use-case:
# Log file analysis

elasticsearch.

# Enter logstash

- Managing events and logs

- Collect data

- Parse data

- Enrich data

- Store data (search and visualizing)

elasticsearch.

# Enter logstash

- Managing events and logs

- Collect data                                  } Input

- Parse data                                    } Filter

- Enrich data

- Store data (search and visualizing) } Output

elasticsearch.

# Data pipeline

- Use a shipper to get your logfiles from all hosts to logstash or a broker (redis, rabbitmq, flume)

- Run data through logstash data pipeline for enrichment

- Store data in elasticsearch

- Use kibana for dashboards and visualisation

elasticsearch.

# Parsing and enrichment?

- Add geo information about an IP

- Parse multi-line exceptions from a java application

- Use grok to have tons of predefined regexes

- Metrics for event throughput information

- HTTP User-Agent extraction

- Enrichment by range values of a field

elasticsearch.

# Use case: Log files

| Logs | Logstash | Store/Search elasticsearch. | Visualize kibana |

elasticsearch.

# Kibana

elasticsearch.

# Kibana

# Kibana

elasticsearch.

# Kibana

# Not-only log files

- Analyse web streams in realtime
  meetup.com RSVP stream
  us gov page visits


- Billing data (payment morale?)


- IRC
  wikipedia changes

elasticsearch.

# Use-case: Analytics

elasticsearch.

# Analytics

- Aggregation of information

- Facets are one dimensional
  Categories/brands/material of all results of this query

- Questions are multidimensional
  Average revenue per category id per day

- Enter Aggregations!

elasticsearch.

# Create knowledge from data

- ## Orders

  How many orders were created every day in the last month?
  How many orders were created per state in the last month?

- ## Money

  What is the average revenue per shopping cart?
  What is the average shopping cart size per order per hour?

- ## Product portfolio

  Take the location of people into account for special offers?
  Analyse page views: Premium or low budget ecommerce site?

elasticsearch.

# Aggregations

```
» curl -X POST 'localhost:9200/orders/order/_search' -d '

{
    "aggs" : {
        "average_order_size" : {
         "avg" : { "field" : "total" }
        }
    }
}
'
```

```
...
  "aggregations" : {
    "average_order_size" : {
      "value" : 658.369
    }
  }
...
```

elasticsearch.

# Aggregations - Filters

```json
{
  "aggs" : {
    "average_order_size_january" : {
      "filter" : {
        "range" : { "created_at" : { "gte" : "2014-01-01", "lt":
"2014-02-01" } } },
      "aggs" : {
        "avg" : { "avg" : { "field" : "total" } }
      }
    }
  }
}
```

```json
...
  "aggregations" : {
    "average_order_size_january" : {
      "doc_count" : 8,
      "avg" : { "value" : 540.89375 }
    }
...
```

elasticsearch.

# Aggregations - per day

```
{
   "aggs": {
     "by_day": {
       "filter": {
         "range": {
           "created_at": {
             "gte": "2014-01-01", "lt": "2014-02-01"
           }
         }
       },
       "aggs": {
         "monthly_filter": {
           "date_histogram": {
             "field": "created_at",
             "interval": "day",
             "format": "yyyy-MM-dd"
           },
           "aggs": {
             "average_order_size": { "avg": { "field": "total" } }
           }
} } } } } }
```

elasticsearch.

# Aggregations - per day

```
...
  "aggregations" : {
    "by_day" : {
      "doc_count" : 8,
      "monthly_filter" : [ {
        "key_as_string" : "2014-01-01",
        "key" : 1388534400000,
        "doc_count" : 136,
        "average_order_size" : {
          "value" : 380.0
        }
      }, {
        "key_as_string" : "2014-01-06",
        "key" : 1388966400000,
        "doc_count" : 256,
        "average_order_size" : {
          "value" : 502.575
        }
      }, {
...
```

elasticsearch.

# Aggregations - per hour

```
{
  "aggs": {
    "by_day": {
      "filter": {
        "range": {
         "created_at": { "gte": "2014-01-01", "lt": "2014-02-01" }
        }
      },
      "aggs": {
        "hourly_filter": {
          "histogram": {
            "interval": 1,
            "script": "doc[\u0027created_at\u0027].date.hourOfDay"
          },
          "aggs": {
            "average_order_size": {
              "avg": { "field": "total" }
            }
          }
        }
      }
} } } }
```

elasticsearch.

# Aggregations - per hour

```
...
  "aggregations" : {
    "by_day" : {
      "doc_count" : 8,
      "hourly_filter" : [ {
        "key" : 11,
        "doc_count" : 1,
        "average_order_size" : {
          "value" : 380.0
        }
      }, {
        "key" : 13,
        "doc_count" : 1,
        "average_order_size" : {
          "value" : 450.15
        }
      }
...
```

elasticsearch.

# Elasticsearch 1.0

elasticsearch.

# Elasticsearch 1.0

- Aggregations

- Snapshot/Restore

- Distributed/scalable percolator

- Cat API
  http://www.elasticsearch.org/blog/introducing-cat-api/

- Federated search: Tribe node

*elasticsearch.*

# Thanks for listening!

elasticsearch.

# Q & A

P.S. We're hiring
http://elasticsearch.com/about/jobs
http://elasticsearch.com/support

Alexander Reelsen
@spinscale
alexander.reelsen@elasticsearch.com

**elasticsearch.**