



FOSDEM 2019



Rated Ranking Evaluator

An open-source approach for Search Quality Evaluation

Alessandro Benedetti, Software Engineer

Andrea Gazzarini, Software Engineer

3rd February 2019





Alessandro Benedetti

- Search Consultant
- R&D Software Engineer
- Master in Computer Science
- Apache Lucene/Solr Enthusiast
- Semantic, NLP, Machine Learning Technologies passionate
- Beach Volleyball Player & Snowboarder





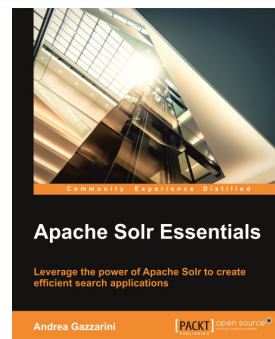
Andrea Gazzarini, “Gazza”



- Software Engineer (1999-)
- “Hermit” Software Engineer (2010-)
- Java & Information Retrieval Passionate
- Apache Qpid (past) Committer
- Husband & Father
- Bass Player



CUMULUSRDF





Search Services

- Open Source Enthusiasts
- Apache Lucene/Solr experts
- Community Contributors
- Active Researchers
- Hot Trends : Learning To Rank, Document Similarity, Search Quality Evaluation, Relevancy Tuning



FAEARCH
INSPIRED FASHION SEARCH

ADOREBEAUTY



BARCLAYS



Alfresco™

Information
Builders

ComplyAdvantage



DeviceBits



NetDimensions

BAR FRIDGES
Cooling the Beers of a Nation Australia



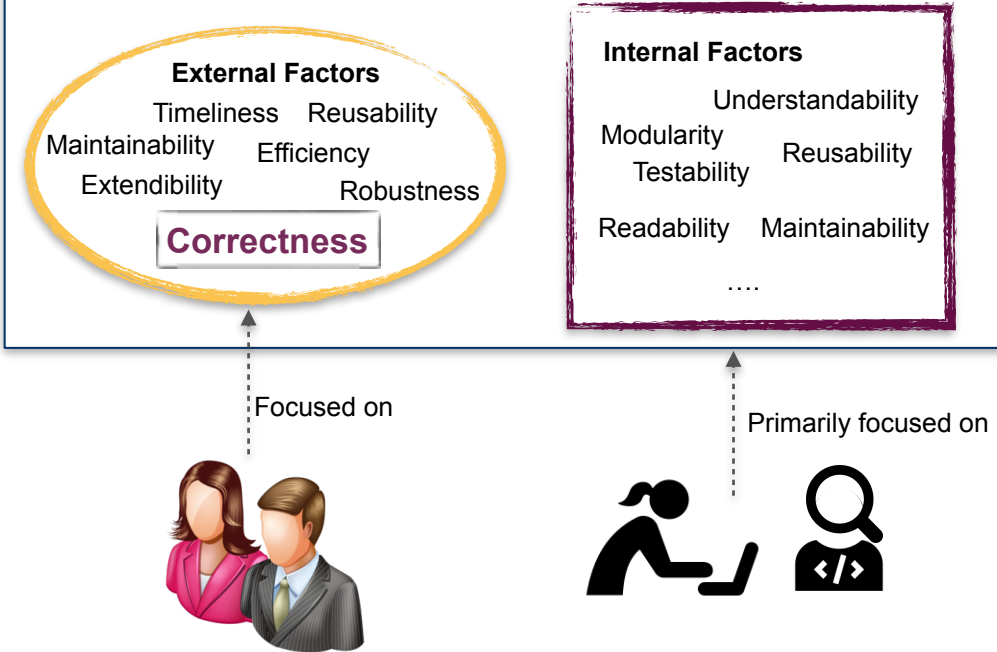


✓ **Search Quality Evaluation**

- **Context overview**
- **Correctness**
- **Evaluation Measures**
- **Rated Ranking Evaluator (RRE)**
- **Demo**
- **Future Works**
- **Q&A**



Search Quality



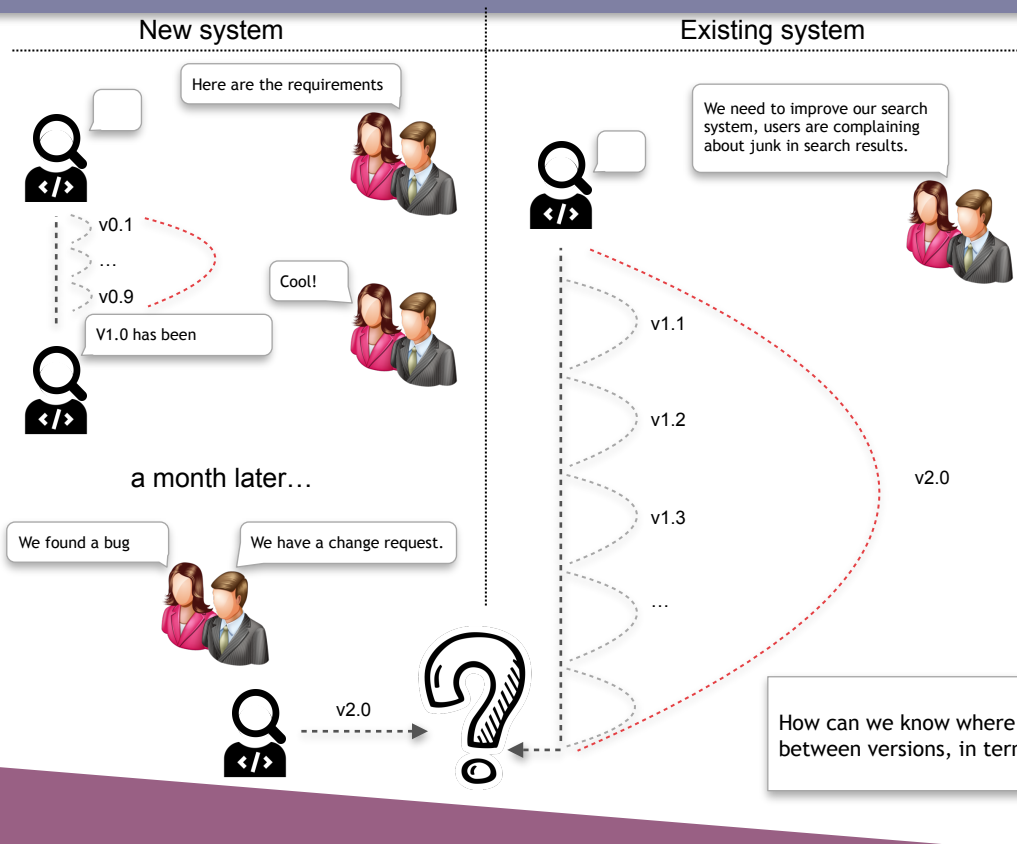
Context Overview

Search engineering is the **production of quality search systems**.

Search quality (and in general software quality) is a huge topic which can be described using **internal** and **external factors**.

In the end, **only external factors matter**, those that can be perceived by **users** and **customers**. But the key for getting **optimal levels** of those external factors are the **internal ones**.

One of the main differences between search and software quality (especially from a **correctness** perspective) is in the **ok / ko judgment**, which is, in general, more “deterministic” in case of software development.



Correctness

Correctness is the ability of a system to perform its exact task, as defined by its specification.

Search domain is critical from this perspective because correctness depends on arbitrary user judgments.

For each internal (gray) and external (red) iteration we need to find **a way to measure the correctness**.

Evaluation measures for an information retrieval system are used to assert how well the search results satisfied the user's query intent.



Evaluation Measures

Offline Measures

Precision
Recall
F-Measure NDCG
Mean Reciprocal Rank
Average Precision
....

Online Measures

Click-through rate
Session abandonment rate
Zero result rate
Session success rate
....

We are mainly focused here

Evaluation Measures

Evaluation measures for an information retrieval system try to formalise how well a search system satisfies its user information needs.

Measures are generally split into two categories: **online** and **offline measures**.

In this context we will focus on **offline measures**.

We will talk about something that can **help a search engineer** during his ordinary day (i.e. in those phases previously called “**internal iterations**”)

We will also see how the **same tool** can be used for a broader usage, like contributing in the **continuous integration pipeline** or even for **delivering value** to **functional stakeholders**.



➤ **Search Quality Evaluation**

✓ **Rated Ranking Evaluator (RRE)**

- ▶ **What is it?**
- ▶ **How does it work?**
- ▶ **Evaluation Process Input & Output**
- ▶ **Challenges**

➤ **Demo**

➤ **Future Works**

➤ **Q&A**



RRE: What is it?



FOSDEM 2019

<https://github.com/SeaseLtd/rated-ranking-evaluator>

SeaseLtd / rated-ranking-evaluator

Unwatch 3 Star 16 Fork 0

Code Issues 10 Pull requests 0 Projects 0 Wiki Insights Settings

Home

Andrea Gazzarini edited this page on Jul 16 · 9 revisions

Rated Ranking Evaluator

The Rated Ranking Evaluator (RRE) is a search quality evaluation tool which, as the name suggests, evaluates the quality of results coming from a search infrastructure.

It is something which helps a Search Engineer in his daily job. Are you a Search Engineer? Are you tuning/implementing/changing/configuring a search infrastructure? Do you want to have something that gives you an evidence about the improvements between changes? So you are in the right place.

There's more: RRE formalises how well a search system satisfies the user information needs, at "technical" level, combining a rich tree-like domain model with several evaluation measures, but also at "functional" level, providing human-readable outputs that could be set as deliverables for

Pages 20

1. What is it?
2. Quick Start
3. Project Structure
4. Evaluation Measures
5. How does it work?
 - 5.1 Domain Model
 - 5.2 What we need to provide
 - 5.3 Where we need to provide
 - 5.4 The Evaluation Process
 - 5.5 The Evaluation Output
6. RRE Server
7. Apache Maven
 - 7.1 Maven Plugin
 - 7.2 Maven Reporting Plugin
 - 7.3 Maven Archetype
8. Search Platform Framework
 - 8.1 Supported platforms
9. FAQ

RRE: What is it?

- A set of **search quality evaluation tools**
- A **search quality evaluation framework**
- **Multi (search) platform**
- Written in **Java**
- It can be used also in **non-Java projects**
- Licensed under **Apache 2.0**
- Open to **contributions**
- **Extremely dynamic!**

RRE: At a glance



FOSDEM 2019

Apache Lucene/Solr
London *meetup*

HIYSTACK
THE SEARCH RELEVANCE CONFERENCE!



FOSDEM 2018

2 Contributors	10 Modules	48950 Lines of Code	2 Months
2 Contributors	10 Modules	67317 Lines of Code	5 Months
4 Contributors	12 Modules	71670 Lines of Code	8 Months





RRE Ecosystem

The picture illustrates the **main modules** composing the RRE ecosystem.

All modules with a **dashed border** are planned for a **future release**.

RRE CLI has a double border because although the **rre-cli** module **hasn't been** developed, you can run **RRE from a command line** using **RRE Maven archetype**, which is part of the current release.

As you can see, the current implementation includes **two** target search platforms: **Apache Solr** and **Elasticsearch**.

The **Search Platform API** module provide a **search platform abstraction** for plugging-in additional search systems.





- ▶ Precision
- ▶ Recall
- ▶ Precision at 1 (P@1)
- ▶ Precision at 2 (P@2)
- ▶ Precision at 3 (P@3)
- ▶ Precision at 10 (P@10)
- ▶ Average Precision (AP)
- ▶ Reciprocal Rank
- ▶ Mean Reciprocal Rank
- ▶ Mean Average Precision (MAP)
- ▶ Normalised Discounted Cumulative Gain (NDCG)
- ▶ **F-Measure** >>>> **Compound Metric**

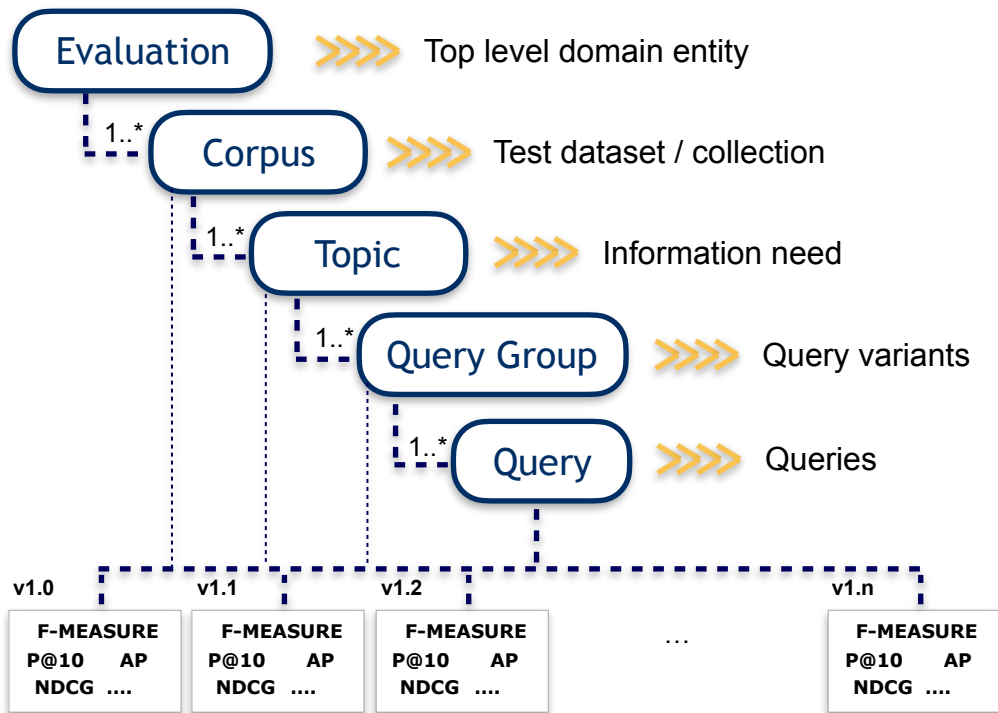
Available Metrics

These are the RRE **built-in metrics** which can be used out of the box.

The most part of them are computed at query level and then aggregated at upper levels.

However, compound metrics (e.g. **MAP**, or **GMAP**) are **not explicitly declared** or defined, because the computation doesn't happen at query level. The result of the aggregation executed on the upper levels will automatically produce these metric.

For example, the **Average Precision** computed for Q1, Q2, Q3, Qn becomes the **Mean Average Precision** at Query Group or Topic levels.



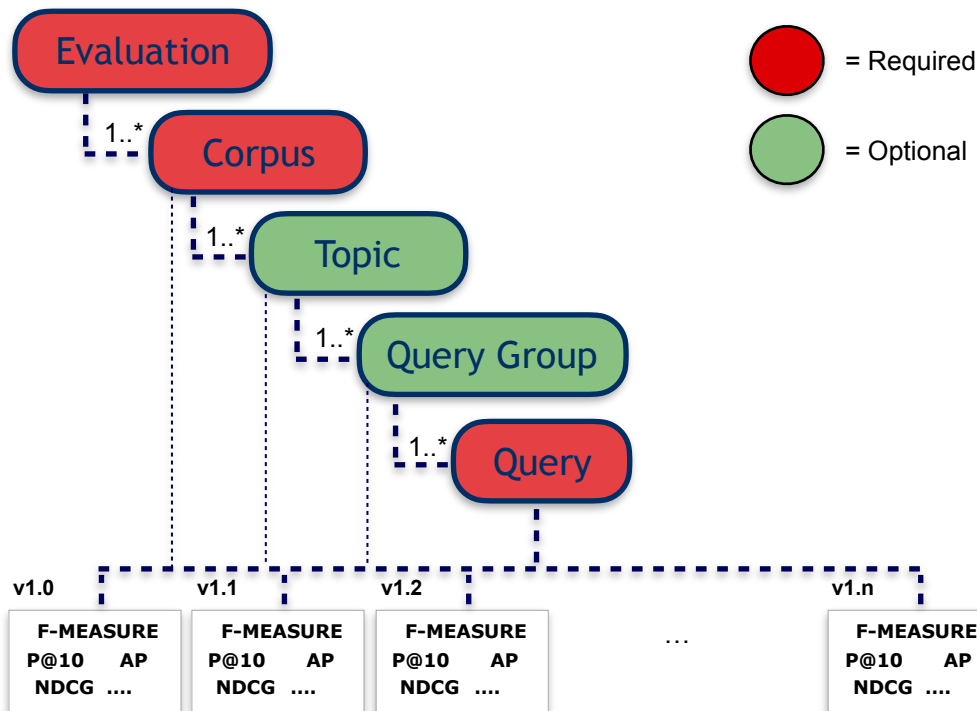
RRE Domain Model

RRE **Domain Model** is organized into a **composite / tree-like** structure where the relationships between entities are **always 1 to many**.

The top level entity is a placeholder representing an **evaluation execution**.

Versioned metrics are computed at **query level** and then reported, using an **aggregation function**, at upper levels.

The benefit of having a composite structure is clear: we can see a **metric value at different levels** (e.g. a query, all queries belonging to a query group, all queries belonging to a topic or at corpus level)



RRE Domain Model

Although the domain model structure is able to **capture complex scenarios**, sometimes we want to model **simpler contexts**.

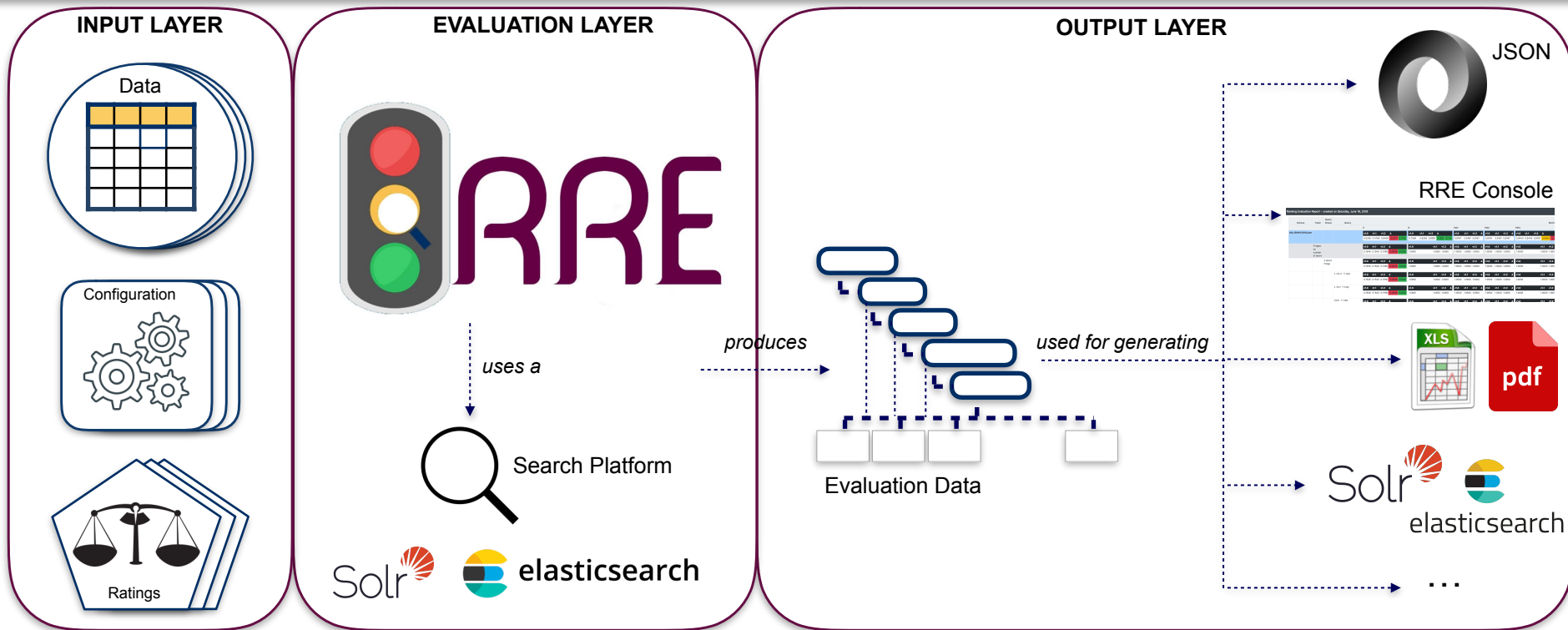
In order to avoid verbose and redundant ratings definitions it's possible to **omit some level**. Specifically we can be in one of the following:

- **only queries**
- **query groups and queries**
- **topics, query groups and queries**

RRE: Evaluation process overview (1/2)



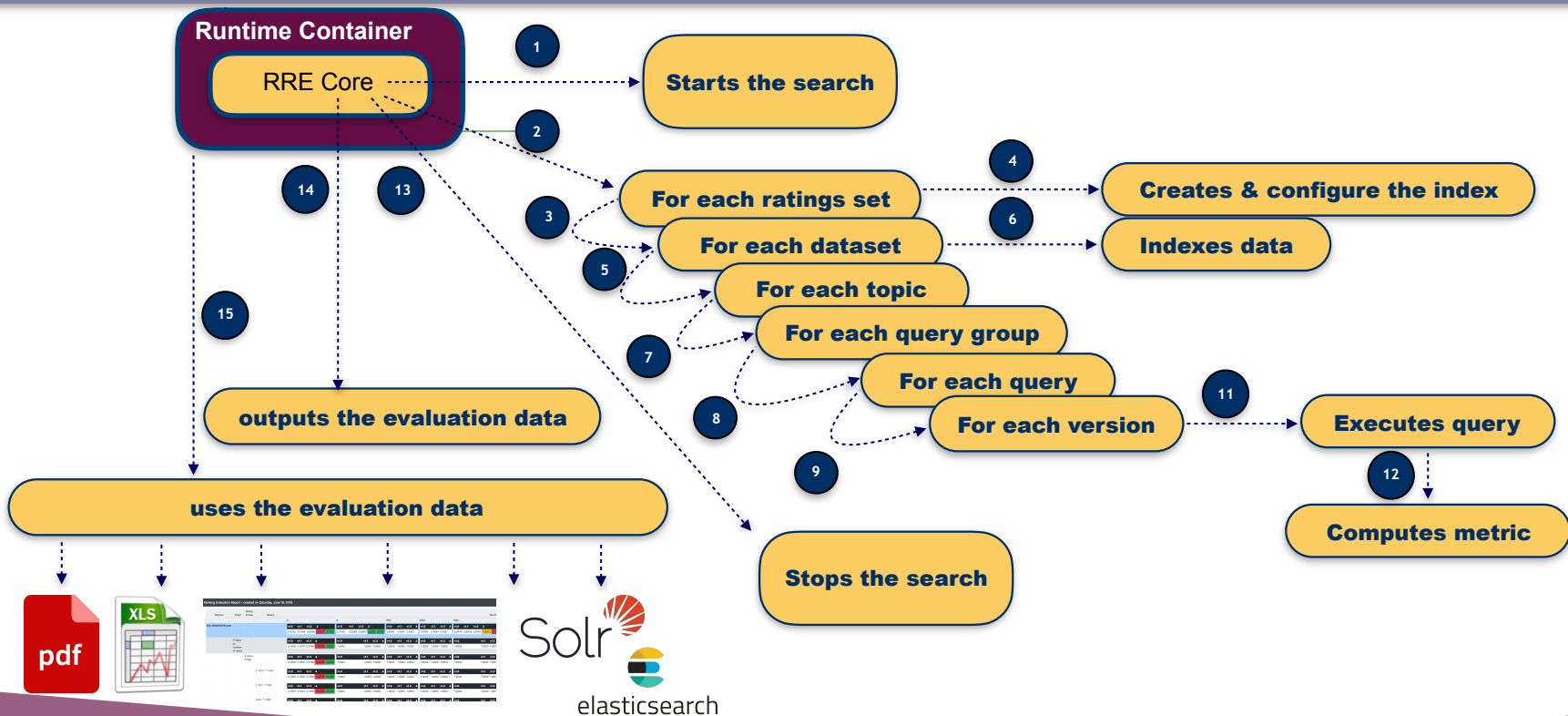
FOSDEM 2019



RRE: Evaluation process overview (2/2)



FOSDEM 2019





```
[
  {
    "id": 1,
    "name": "Fender Jazz Bass"
  },
  {
    "id": 2,
    "name": "Fender Precision Bass"
  },
  {
    "id": 3,
    "name": "Warwick Corvette"
  },
  {
    "id": 4,
    "name": "Warwick Thumb"
  }
]
```



```
{
  "index" : { "_index": "dataset1", "_type": "doc", "_id": "1" }
  "name": "Fender Jazz Bass"
}
{
  "index" : { "_index": "dataset1", "_type": "doc", "_id": "2" }
  "name": "Fender Precision Bass"
}
{
  "index" : { "_index": "dataset1", "_type": "doc", "_id": "3" }
  "name": "Warwick Corvette"
}
{
  "index" : { "_index": "dataset1", "_type": "doc", "_id": "4" }
  "name": "Warwick Thumb"
}
```

Corpora

An evaluation execution can involve **more than one datasets** targeting a given **search platform**.

A dataset consists consists of **representative** domain data; although a compressed dataset can be provided, generally it has a small/medium size.

Within RRE, **corpus**, **dataset**, **collection** are synonyms.

Datasets must be **located** under a **configurable folder**. Each dataset is then referenced in one or more **ratings file**.



- ▼ 1.1
 - index-shape.json
- ▼ 1.2
 - index-shape.json
- ▼ 1.3
 - en_protected_keywords.txt
 - index-shape.json



elasticsearch

- ▼ v1.0
 - products
 - conf
 - lang
 - dih-config.xml
 - import.properties
 - mapping-FoldToASCII.txt
 - schema.xml
 - solrconfig.xml
 - units.json
 - solr.xml
- ▼ v1.1
 - products
 - conf
 - lang
 - dih-config.xml
 - mapping-FoldToASCII.txt
 - schema.xml
 - solrconfig.xml
 - units.json
 - solr.xml
- ▼ v1.2
 - products

Configuration Sets

The search platform **configuration evolves over time** (e.g. change requests, enhancements, bugs)

RRE encourages an **incremental approach** for managing the configuration instances. Even for internal or small iterations, each time we make a relevant change to the current configuration, it's better to **clone it and move forward** with a new version.

In this way we'll end up having the **historical progression** of our system, and RRE will be able to **make comparisons**.

The evaluation process allows you to define **inclusion / exclusion** rules (i.e. include only version 1.0 and 2.0)





```
"queries": [
  {
    "template": "only_q.json",
    "placeholders": {
      "$query": "fender"
    }
  },
  {
    "template": "filter_by_language.json",
    "placeholders": {
      "$query": "Fender",
      "$lang": "eng"
    }
  }
]

"topics": [
  {
    "description": "Out of warranty",
    "query_groups": [
      {
        "name": "Out of warranty variants",
        "description": "The group tests several v",
        "template": "query-shape-${version}.json",
        "queries": [
          {
            "placeholders": {
              "$query##": "oow"
            }
          },
          {
            "placeholders": {
              "$query##": "0ow"
            }
          }
        ]
      }
    ]
  }
]
```

only_q.json

```
{
  "q" : "$query"
}
```

filter_by_language.json

```
{
  "q" : "$query",
  "fq" : "language:$lang"
}
```

function_score: {

```
  "query": {
    "bool": {
      "should": [
        {
          "constant_score": {
            "filter": {
              "match": {
                "classification_language_1": "#$query##"
              }
            },
            "boost": 3
          }
        }
      ]
    }
  },
  "constant_score": {
    "filter": {
      "match": {
        "classification_language_1": "#$query##"
      }
    }
  }
}
```

Query templates

For each **query** or **query group**) it's possible to define a **template**, which is a kind of **query shape** containing one or more **placeholders**.

Then, in the ratings file you can reference one of those defined templates and you can provide a **value for each placeholder**.

Templates have been introduced in order to:

- allow a **common query management** between search platforms
- define **complex queries**
- define **runtime parameters** that cannot be statically determined (e.g. filters)



```
{
  "index" : "core1",
  "corpora_file": "electric_basses.json",
  "id_field": "id",
  "topics" : [
    {
      "description": "Fender basses",
      "query_groups": [
        {
          "name": "Brand search",
          "description": "The group tests several searches on the Fender brand",
          "queries": [
            {
              "template": "only_q.json",
              "placeholders": {
                "$query": "fender"
              }
            },
            {
              "template": "only_q.json",
              "placeholders": {
                "$query": "Fender"
              }
            }
          ]
        },
        {
          "name": "Jazz bass search",
          "description": "Several searches on a given model (Jazz bass)",
          "queries": [
            {
              "template": "only_q.json",
              "placeholders": {
                "$query": "Jazz bass"
              }
            }
          ]
        }
      ]
    }
  ]
}
```

OR

```
{
  "relevant_documents": {
    "1": {
      "gain": 3
    },
    "2": {
      "gain": 3
    }
  }
}
```

Ratings

Ratings files associate the RRE domain model **entities** with **relevance judgments**. A ratings file provides the association **between queries and relevant documents**.

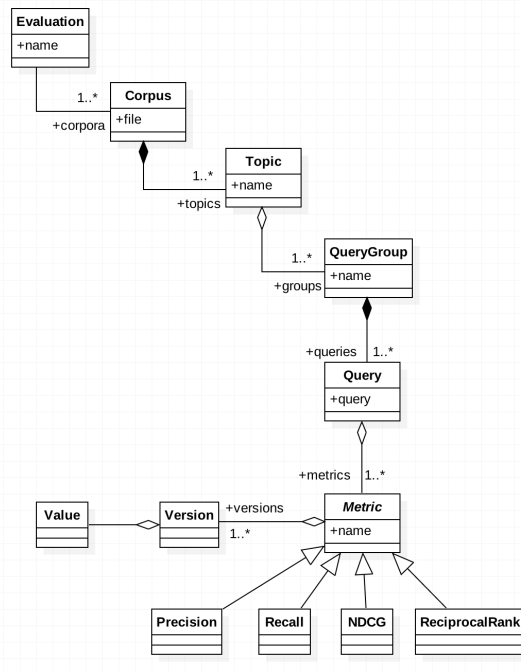
There must be **at least one** ratings file (otherwise no evaluation happens). Usually there's a 1:1 relationship between a rating file and a dataset.

Judgments, the most important part of this file, consist of a **list of all relevant documents** for a **query group**.

Each listed document has a corresponding **“gain”** which is the **relevancy judgment** we want to assign to that document.



```
{
  "name": "Ranking Evaluation Report - created on Friday, June 15, 2018",
  "metrics": { ... },
  "corpora": [ {
    "name": "bfa_15MAY2018.json",
    "metrics": { ... },
    "topics": [ {
      "name": "Fridges by number of doors",
      "metrics": { ... },
      "query-groups": [ {
        "name": "2 doors fridge",
        "metrics": { ... },
        "query-evaluations": [ {
          "metrics": { ... },
          "query": "2 doors fridge",
          "results": { ... }
        }
      ]
    }
  ],
  "metrics": {
    "p": {
      "name": "P",
      "versions": {
        "v1.0": {
          "value": "0.1546"
        },
        "v1.1": {
          "value": "0.1539"
        },
        "v1.2": {
          "value": "0.1798"
        }
      }
    },
    "trend": [ -0.0007, 0.0259 ]
  },
  "r": {
    "name": "R",
    "versions": {
      "v1.0": {
        "value": "1.0000"
      },
      "v1.1": {
        "value": "1.0000"
      },
      "v1.2": {
        "value": "1.0000"
      }
    },
    "trend": [ 0.0000, 0.0000 ]
  },
  "ap": {
    "name": "AP",
    "versions": {
      "v1.0": {
        "value": "1.00000000"
      }
    }
  }
}
```



Evaluation output

The **RRE Core** itself is a **library**, so it outputs its result as a **Plain Java object** that must be programmatically used.

However when wrapped within a **runtime container**, like the **Maven Plugin**, the evaluation object tree is marshalled in **JSON format**.

Being interoperable, the JSON format can be used by some other component for producing a different **kind of output**.

An example of such usage is the **RRE Apache Maven Reporting Plugin** which can

- output a **spreadsheet**
- send the evaluation data to a running **RRE Server**



Workbook

The RRE domain model (topics, groups and queries) is on the **left** and each metric (on the right section) has a **value** for each **version / entity** pair.

In case the evaluation process includes multiple datasets, there will be a **spreadsheet** for each of them.

This output format is useful when

- you want to have (or maintain somewhere) a **snapshot** about how the system performed in a given moment
- the comparison includes a **lot of versions**
- you want to include **all available metrics**

Topic	Query Group	Query	Metric	NDCG@10										DELTA		0.8	0.9
				0.8	0.9	1.1	1.2	1.3	1.4	1.5	0.8772	0.5644	0.0	0.0	0.0		
1				0.4372	0.36	0.9244	0.9244	0.9244	0.9244	0.8844	-0.0772	0.5644	0.0	0.0	0.0	0.278	0.36
5	SIM issues			0	1	1	1	1	1	1	1	1	0	0	0	0	1
6		The group tests several variants of the topic above.		0	1	1	1	1	1	1	1	1	0	0	0	0	1
7		Sim not readable		0	1	1	1	1	1	1	1	1	0	0	0	0	1
8		SIM unreadable		0	1	1	1	1	1	1	1	1	0	0	0	0	1
9	Long distance calls			0	1	1	1	1	1	1	1	1	0	0	0	0	1
10		The group tests several variants of the topic above.		0	1	1	1	1	1	1	1	1	0	0	0	0	1
11		Long distance calls		0	1	1	1	1	1	1	1	1	0	0	0	0	1
12		Long distance models		0	1	1	1	1	1	1	1	1	0	0	0	0	1
13	Samsung Galaxy S7 manual			0.5134	0	1	1	1	1	1	-0.5134	1	0	0	0	0.5334	0
14		The group tests several variants of the topic above.		0.5134	0	1	1	1	1	1	-0.5134	1	0	0	0	0.5334	0
15		Tutorial for S7		0.68	0	1	1	1	1	1	-0.68	1	0	0	0	0.6	0
16		S7 manuals		0.43	0	1	1	1	1	1	-0.43	1	0	0	0	0.5	0
17		S7 Manuals		0.43	0	1	1	1	1	1	-0.43	1	0	0	0	0.5	0
18	Unlock S7			1	0	1	1	1	1	1	-1	1	0	0	0	0.1	0
19		The group tests several variants of the topic above.		1	0	1	1	1	1	1	-1	1	0	0	0	0.1	0
20		Unlock S7		1	0	1	1	1	1	1	-1	1	0	0	0	0.1	0
21		Unlock Galaxy S7		1	0	1	1	1	1	1	-1	1	0	0	0	0.1	0
22		Unlock Samsung S7		1	0	1	1	1	1	1	-1	1	0	0	0	0.1	0
23		Unlock Galaxy Samsung S7		1	0	1	1	1	1	1	-1	1	0	0	0	0.1	0
24	Unlock (Specific for S7 Edge)			0.75	0	1.25	1.25	1.25	1.25	1.25	-0.75	1.25	0	0	0	0.1	0
25		The group tests several variants of the topic above.		0.75	0	1.25	1.25	1.25	1.25	1.25	-0.75	1.25	0	0	0	0.1	0
26		Unlock S7 edge		0.75	0	1.25	1.25	1.25	1.25	1.25	-0.75	1.25	0	0	0	0.1	0
27		s7 edge unlock		0.75	0	1.25	1.25	1.25	1.25	1.25	-0.75	1.25	0	0	0	0.1	0
28	Screen Frozen			1	0	1	1	1	1	1	-1	1	0	0	0	1	0
29		The group tests several variants of the topic above.		1	0	1	1	1	1	1	-1	1	0	0	0	1	0
30		screen frozen can't power off		1	0	1	1	1	1	1	-1	1	0	0	0	1	0
31		screen frozen		1	0	1	1	1	1	1	-1	1	0	0	0	1	0
32	Cancel Auto-refill			0	1	1	1	1	1	1	1	0	0	0	0	0	1
33		The group tests several variants of the topic above.		0	1	1	1	1	1	1	1	0	0	0	0	0	1
34		Cancel auto refill		0	1	1	1	1	1	1	1	0	0	0	0	0	1
35		cancel autorefill		0	1	1	1	1	1	1	1	0	0	0	0	0	1
36		cancel auto-refill		0	1	1	1	1	1	1	1	0	0	0	0	0	1
37		how to cancel auto-refill		0	1	1	1	1	1	1	1	0	0	0	0	0	1
38		how to cancel autorefill		0	1	1	1	1	1	1	1	0	0	0	0	0	1
39	Text messaging with VIVO 5			0.7467	0	1	1	1	1	1	-0.7467	1	0	0	0	0.3	0
40		The group tests several variants of the topic above.		0.7467	0	1	1	1	1	1	-0.7467	1	0	0	0	0.3	0
41		Text messaging with VIVO 5		0.93	0	1	1	1	1	1	-0.93	1	0	0	0	0.4	0
42		text message VIVO 5		1	0	1	1	1	1	1	-1	1	0	0	0	0.4	0
43		sms vivo 5		0.31	0	1	1	1	1	1	-0.31	1	0	0	0	0.1	0
44	Text messaging (send) - VIVO			0.75	0	1	1	1	1	1	-0.75	1	0	0	0	0.1	0
45		The group tests several variants of the topic above.		0.75	0	1	1	1	1	1	-0.75	1	0	0	0	0.1	0
46		sendin a text message using		0.4	0	1	1	1	1	1	-0.4	1	0	0	0	0.1	0



Ranking Evaluation Report - created on Wednesday, June 20, 2018

Corpus	Topic	Query Group	Query	P				R				RR@10				AP				Metrics NDCG@10			
				v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ
bfa_15MAY2018.json				0.5230	0.4759	0.6459	0.2471	0.7026	0.8293	0.9397	0.2371	0.9443	0.9443	0.9376	0.0000	0.8696	0.7895	0.8962	0.1067	0.8801	0.8734	0.9050	0.0249
	Fridges by number of doors			0.1548	0.1540	0.1798	0.0200	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000
		2 doors fridge		0.1548	0.1540	0.1798	0.0200	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000
		2 doors fridge		0.1548	0.1540	0.1798	0.0200	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000
		2 door fridge		0.1548	0.1540	0.1798	0.0200	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000
		2door fridge		0.1548	0.1540	0.1798	0.0200	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000
		2doors fridge		0.1548	0.1540	0.1798	0.0200	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000	1.0000	1.0000	1.0000	0.0000
	Beer fridge quiet indoor			0.0201	0.0149	0.0297	0.0552	0.5000	0.6000	1.0000	0.5000	0.3301	0.3301	0.2500	0.0800	0.1827	0.1827	0.6022	0.4200	0.3401	0.3401	0.5800	0.2399
		Beer fridge quiet indoor		0.0201	0.0149	0.0297	0.0552	0.5000	0.6000	1.0000	0.5000	0.3301	0.3301	0.2500	0.0800	0.1827	0.1827	0.6022	0.4200	0.3401	0.3401	0.5800	0.2399
		Beer fridge quiet indoor		0.0201	0.0149	0.0297	0.0552	0.5000	0.6000	1.0000	0.5000	0.3301	0.3301	0.2500	0.0800	0.1827	0.1827	0.6022	0.4200	0.3401	0.3401	0.5800	0.2399
	Coloured mini fridges			0.3182	0.0401	1.0000	0.2781	0.2188	0.2188	0.8750	0.6562	1.0000	1.0000	1.0000	0.0000	0.1797	0.1485	0.8750	0.7265	0.6100	0.6700	1.0000	0.3300
	black mini fridges			0.3182	0.0401	1.0000	0.2781	0.2188	0.2188	0.8750	0.6562	1.0000	1.0000	1.0000	0.0000	0.1797	0.1485	0.8750	0.7265	0.6100	0.6700	1.0000	0.3300

RRE Server

The RRE console is a **SpringBoot/AngularJS** application which shows real-time information about **evaluation results**.

Each time a build happens, the **RRE reporting plugin** sends the evaluation result to a **RESTful endpoint** provided by RRE Server.

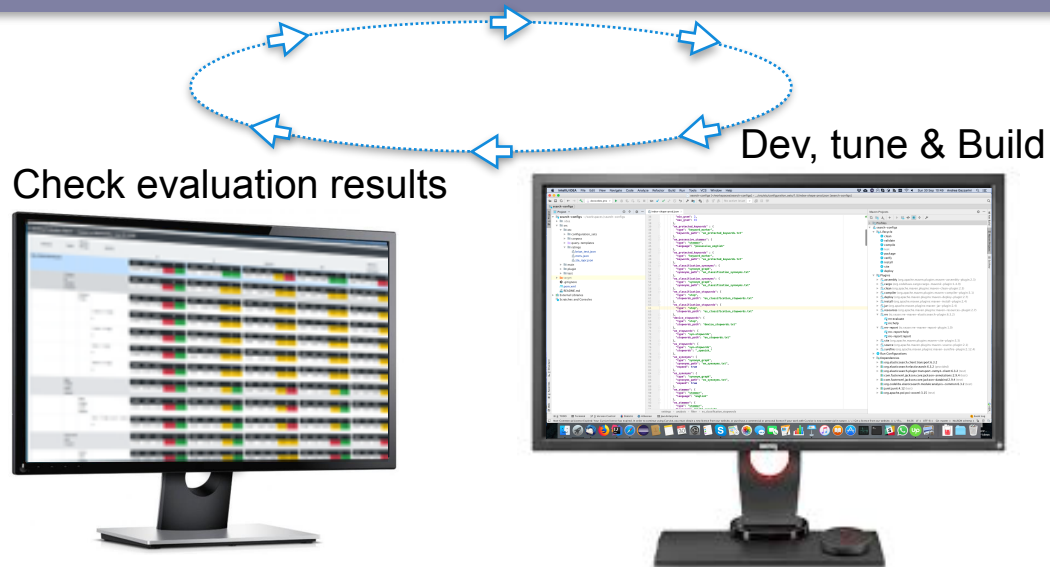
The received data **immediately updates** the web dashboard with fresh data.

Useful during the **development / tuning phase iterations** (you don't have to open again and again the excel report)

RRE: Iterative development & tuning



FOSDEM 2019



We are thinking about how to fill a third monitor





“I think if we could create a simplified **pass/fail** report for the business team, that would be ideal. So they could understand the tradeoffs of the new search.”



Do I have to write all judgments manually??



“Many search engines process the user query heavily before it's submitted to the search engine in whatever DSL is required, and if you don't retain some idea of the original query in the system how can you” relate the test results back to user behaviour?

How can I use RRE if I have a custom search platform?

Can I persist the evaluation data?

Java is not in my stack



- Search Quality Evaluation
- Rated Ranking Evaluator
- ✓ **Demo**
- Future Works
- Q&A

RRE DEMO: Github Repository



FOSDEM 2019

<https://github.com/SeaseLtd/rre-fosdem-walkthrough>

The screenshot shows the Github repository page for 'SeaseLtd / rre-fosdem-walkthrough'. The repository has 1 commit, 1 branch, 0 releases, and 1 contributor. The 'master' branch is selected. A table of recent commits is shown, with the latest commit by 'agazzarini' merging branch 'master' of the repository, dated 11 days ago. The commit message is 'Merge branch 'master' of https://github.com/SeaseLtd/rre-fosdem-walk...'. The commit details show files: 'src/etc' (README files, 11 days ago), '.gitignore' (Examples from Ch 3, 13 days ago), 'README.md' (Update README.md, 11 days ago), and 'pom.xml' (README files, 11 days ago). The README.md file is selected, showing the RRE logo and the text 'RRE Walkthrough Demo Repository'.

File	Commit Message	Time
src/etc	README files	11 days ago
.gitignore	Examples from Ch 3	13 days ago
README.md	Update README.md	11 days ago
pom.xml	README files	11 days ago

RRE Walkthrough Demo Repository

Example Project

- A sample RRE-enabled project
- No Java code, only configuration
- Search Platform: Elasticsearch 6.3.2
- Seven example iterations
- Index shapes & queries from Relevant Search [1]
- Dataset: TMBD (extract)

[1] <https://www.manning.com/books/relevant-search>



RRE DEMO: Let's start!



FOSDEM 2019

Project Skeleton

- Project skeleton with example files
- Generated using RRE Maven Archetype
- Can be executed in the command line
- Can be imported in an IDE (e.g. IntelliJ, Eclipse)
- Prerequisites: Java \geq 8, Apache Maven 3.x
- Detailed steps in the RRE Wiki



RRE artefacts are not yet available in central Maven Repository. As consequence of that, Sease Repository needs to be configured in your Maven settings. The RRE Wiki contains detailed info about that.



github.com/SeaseLtd/rated-ranking-evaluator/wiki/Quick%20Start

SeaseLtd / [rated-ranking-evaluator](#)

Unwatch 7 Unstar 51 Fork 3

Code Issues 12 Pull requests 0 Projects 0 Wiki Insights Settings

Quick Start

Andrea Gazzarini edited this page on Sep 28, 2018 · 6 revisions

The quickest way to start with RRE is through the available Maven plugins. In this context, you can be in one of the following scenarios:

- **An existing project which already uses Java and Apache Maven:** perfect, all you need to do is to configure the RRE plugin in the pom.xml
- **A new project where you will be using Java and Apache Maven:** still perfect, the RRE Maven archetype creates a project skeleton with all what you need
- **Existing or new Java project but no Apache Maven:** other build management tools are planned for future releases, so the Maven archetype is the right path to follow. It will create a project with all configuration needed for running the evaluation process; you don't need to put any code inside that project, your code will live somewhere else, with some other tool (e.g. Gradle)
- **Existing or new non-Java project:** the Maven archetype is still the right choice. With a minimal set of pre-requirements (i.e. Java and Apache Maven), you can create an RRE-enabled module for running evaluation builds.

So at the end, in order to have a first taste of RRE, we need to create a project from scratch (by means of the RRE Maven Archetype) or to configure an existing project. Let's detail both alternatives.

New Project

Pages 20

1. What is it?
2. Quick Start
3. Project Structure
4. Evaluation Measures
5. How does it work?
 - 5.1 Domain Model
 - 5.2 What we need to provide
 - 5.3 Where we need to provide
 - 5.4 The Evaluation Process
 - 5.5 The Evaluation Output
6. RRE Server
7. Apache Maven
 - 7.1 Maven Plugin
 - 7.2 Maven Reporting Plugin
 - 7.3 Maven Archetype
8. Search Platform Framework
 - 8.1 Supported platforms
9. FAQ

Clone this wiki locally

<https://github.com/Sease>



```
> mvn archetype:generate \  
-PSease \  
-B \  
-DarchetypeGroupId=io.sease \  
-DarchetypeArtifactId=rre-maven-elasticsearch-archetype \  
-DarchetypeVersion=1.0 \  
-DgroupId=com.yourcompany.rre \  
-DartifactId=search-quality-evaluation \  
-Dversion=1.0 \  
-DesVersion=6.3.2
```



Apache Solr folks will use the following properties:

```
...  
-DarchetypeArtifactId=rre-maven-solr-archetype \  
...  
-DsolrVersion=7.4.0
```

Maven Archetype

- Prerequisites: **Java** >= 8, **Apache Maven** 3.x
- **-B** stands for batch mode. Without that parameter the skeleton build procedure will be **interactive**
- **groupId** and **artifactId** can be whatever
- **archetypeArtifactId**, **esVersion/solrVersion** need to be configured according with your search infrastructure



- RRE enabled with selected Search Platform
- Two sample iterations (v1.0, v1.1)
- Each folder contains
 - an example of the expected content
 - a README.md



Starter Project

- In the pom.xml you can customise:
 - the evaluation metrics
 - the project folders structure
 - the search platform version (e.g. 6.3.2)



You could even change the search platform (e.g. from Elasticsearch to Apache Solr, in the example). Note that without restarting from scratch, changing the search platform type on an existing project would invalidate the example files.

```
index-shape.json x search-quality-evaluation x
1 <?xml version="1.0" encoding="UTF-8" ?>
2 <project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
3       xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org/xsd/maven-4.0.0.xsd">
4   <modelVersion>4.0.0</modelVersion>
5   <groupId>com.yourcompany.rre</groupId>
6   <artifactId>search-quality-evaluation</artifactId>
7   <version>1.0</version>
8   <packaging>pom</packaging>
9   <build>
10     <plugins>
11       <plugin>
12         <groupId>io.sease</groupId>
13         <artifactId>rre-maven-elasticsearch-plugin</artifactId>
14         <version>6.3.2</version>
15         <!-- the configuration below is provided just for example, as it perfectly matches default values -->
16         <configuration>
17           <configurations-folder>src/etc/configuration_sets</configurations-folder>
18           <corpora-folder>src/etc/corpora</corpora-folder>
19           <ratings-folder>src/etc/ratings</ratings-folder>
20           <templates-folder>src/etc/templates</templates-folder>
21           <fields>*,score</fields>
22           <metrics>
23             <param>io.sease.rre.core.domain.metrics.impl.Precision</param>
24             <param>io.sease.rre.core.domain.metrics.impl.Recall</param>
25             <param>io.sease.rre.core.domain.metrics.impl.ReciprocalRanks</param>
26             <param>io.sease.rre.core.domain.metrics.impl.AveragePrecision</param>
27             <param>io.sease.rre.core.domain.metrics.impl.NDCGAtTen</param>
28             <param>io.sease.rre.core.domain.metrics.impl.PrecisionAtOne</param>
29             <param>io.sease.rre.core.domain.metrics.impl.PrecisionAtTwo</param>
30             <param>io.sease.rre.core.domain.metrics.impl.PrecisionAtThree</param>
31             <param>io.sease.rre.core.domain.metrics.impl.PrecisionAtTen</param>
32           </metrics>
33         </configuration>
34         <executions>
35           <execution>
36             <id>search-quality-evaluation</id>
37             <phase>package</phase>
38             <goals>
39               <goal>evaluate</goal>
40             </goals>
41           </execution>
42         </executions>
43       </plugin>
44     </plugins>
45   </build>
46 </project>
```



https://github.com/SeaseLtd/rre-fosdem-walkthrough/tree/master/src/etc/configuration_sets/v1.0

```
1  {
2    "settings": {
3      "number_of_shards": 1,
4      "number_of_replicas": 0
5    },
6    "mappings": {
7      "doc": {
8
9      }
10   }
11 }
```

```
    {
      "index": "tbmd",
      "corpora_file": "tbmd.bulk",
      "id_field": "_id",
      "topics": [
        {
          "description" : "Space Jam",
          "query_groups": [
            {
              "name": "See \"Relevant Search\", Chapter 3",
              "queries": [
                {
                  "template": "multi_match_on_title.json",
                  "placeholders": {
                    "$query": "basketball with cartoon aliens"
                  }
                }
              ]
            }
          ],
          "relevant_documents": {
            "2300": {
              "gain": 3
            }
          }
        }
      ]
    }
  }
```

<https://github.com/SeaseLtd/rre-fosdem-walkthrough/tree/master/src/etc/ratings>

tbmd.json

<https://github.com/SeaseLtd/rre-fosdem-walkthrough/tree/master/src/etc/corpora>

What we need

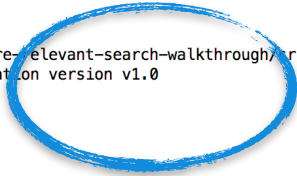
- At least the first iteration (version)
- Corpus
- Ratings
- Query templates



You can directly load in your IDE the demo repository or copy and paste gradually the required data from there.



```
[INFO] Building rre-fosdem-walkthrough 1.0
[INFO] -----
[INFO] --- rre-maven-elasticsearch-plugin:6.3.2:evaluate (search-quality-evaluation) @ rre-fosdem-walkt
[INFO] RRE: New evaluation session is starting...
[INFO] RRE: Search Platform in use: Elasticsearch
[INFO] RRE: Starting Elasticsearch...
[INFO] RRE: Elasticsearch Search Platform successfully started.
[INFO] RRE: found 1 ratings sets.
[INFO] RRE: Ratings Set processing starts
[INFO] RRE: *****
[INFO] RRE: Index name => tbmd
[INFO] RRE: ID Field name => _id
[INFO] RRE: Test Collection => /Users/agazzarini/workspaces/fosdem/rre-relevant-search-walkthrough/r
[INFO] RRE: Loading the Test Collection into Elasticsearch, configuration version v1.0
[INFO] RRE: Elasticsearch has been correctly loaded.
[INFO] RRE: target versions are v1.0
[INFO] TOPIC: Space Jam
[INFO] QUERY GROUP: See "Relevant Search", Chapter 3
[INFO] QUERY: basketball with cartoon aliens
[INFO] TOPIC: Star Trek
[INFO] QUERY GROUP: Starring
[INFO] QUERY: patrick stewart
[INFO] QUERY GROUP: Starring and Title
[INFO] QUERY: star trek patrick stewart
[INFO] QUERY GROUP: Names and Title
[INFO] QUERY: star trek patrick stewart william shatner
[INFO] RRE: Elasticsearch Search Platform shutdown procedure executed.
[INFO] --- rre-maven-report-plugin:1.0:report (search-quality-evaluation-reporting) @ rre-fosdem-walkt
[INFO] Evaluation data has been correctly sent to RRE Server located at http://127.0.0.1:8080
[INFO] --- maven-install-plugin:2.4:install (default-install) @ rre-fosdem-walkthrough ---
[INFO] Installing /Users/agazzarini/workspaces/fosdem/rre-relevant-search-walkthrough/pom.xml to /User
[INFO] -----
[INFO] BUILD SUCCESS
```



Iteration #1

- Default configuration (empty config)
- See the execution plan in the logs
- Two different outputs:
 - Excel
 - RRE Server



The output type can be configured in the pom.xml

RRE DEMO: Iteration #1 (2/3)



FOSDEM 2019

Ranking Evaluation Report - created on Sunday, January 27, 2019

Corpus	Topic	Query Group	Query	Metrics							
				P@1		P@10		AP		NDCG@10	
				v1.0	Δ	v1.0	Δ	v1.0	Δ	v1.0	Δ
tmbd.bulk				0.5000	0.5000	0.3750	0.3750	0.2853	0.2853	0.7126	0.7126
	Space Jam			0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		See "Relevant Search", Chapter 3		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
			basketball with cartoon aliens	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
	Star Trek			0.6667	0.6667	0.5000	0.5000	0.3804	0.3804	0.9500	0.9500
		Starring		0.0000	0.0000	0.2001	0.2001	0.0430	0.0430	0.1601	0.1601
			patrick stewart	0.0000	0.0000	0.2001	0.2001	0.0430	0.0430	0.1601	0.1601
		Starring and Title		1.0000	1.0000	0.7000	0.7000	0.5686	0.5686	0.8001	0.8001
			star trek patrick stewart								

RRE DEMO: Iteration #1 (3/3)



FOSDEM 2019

🔍 Ranking Evaluation Report - created on Sunday, January 27, 2019

[illegible]

OMG! Only red items :(





Iteration #2 (v1.1)

- Index shape contains some specific settings
- Default **english** analyzer for **title** and **overview** fields
- mmmm...everything is still red
- ...however results are different

```
[INFO] --- rre-maven-elasticsearch-plugin:6.3.2:evaluate (search-quality-evaluation) @ rre-fosdem-walkthrough ---
[INFO] RRE: New evaluation session is starting...
[INFO] RRE: Search Platform in use: Elasticsearch
[INFO] RRE: Starting Elasticsearch...
[INFO] RRE: Elasticsearch Search Platform successfully started.
[INFO] RRE: found 1 ratings sets.
[INFO] RRE: Ratings Set processing starts
[INFO] *****
[INFO] RRE: Index name => tmbd
[INFO] RRE: ID Field name => _id
[INFO] RRE: Test Collection => /Users/agazzarini/workspaces/fosdem/rre-relevant-search-walkthrough/src/etc/corpora/tmbd.bulk
[INFO] RRE: Loading the Test Collection into Elasticsearch, configuration version v1.0
[INFO] RRE: Loading the Test Collection into Elasticsearch, configuration version v1.1
[INFO] RRE: Elasticsearch has been correctly loaded.
[INFO] RRE: target versions are v1.0,v1.1
[INFO] TOPIC: Space Jam
[INFO] QUERY GROUP: See "Relevant Search", Chapter 3
[INFO] QUERY: basketball with cartoon aliens
[INFO] TOPIC: Star Trek
[INFO] QUERY GROUP: Starring
[INFO] QUERY: patrick stewart
[INFO] QUERY GROUP: Starring and Title
[INFO] QUERY: star trek patrick stewart
[INFO] QUERY GROUP: Names and Title
[INFO] QUERY: star trek patrick stewart william shatner
[INFO] RRE: Elasticsearch Search Platform shutdown procedure executed.
[INFO] --- rre-maven-report-plugin:1.0:report (search-quality-evaluation-reporting) @ rre-fosdem-walkthrough ---
[INFO] Evaluation data has been correctly sent to RRE Server located at http://127.0.0.1:8080
```

RRE DEMO: Iteration #1 (2/2)



FOSDEM 2019

Ranking Evaluation Report - created on Sunday, January 27, 2019

Corpus	Topic	Query Group	Query	Metrics											
				P@1			P@10			AP			NDCG@10		
				v1.0	v1.1	Δ	v1.0	v1.1	Δ	v1.0	v1.1	Δ	v1.0	v1.1	Δ
tmbd.bulk				0.5000	0.0000	-0.5000	0.3750	0.3750	0.0000	0.2853	0.2448	-0.0405	0.7126	0.6725	-0.0401
Space Jam				0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		See "Relevant Search", Chapter 3		0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
			basketball with cartoon aliens	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Ver. v1.0 Hits: 1197	#1 title Aliens id 679	#2 title The Basketball Diaries id 10474	#3 title Cowboys & Aliens id 49849	#4 title Monsters vs Aliens id 15512	#5 title Aliens vs Predator: Requiem id 440	#6 title Aliens in the Attic id 20856	#7 title Dances with Wolves id 581	#8 title Friends with Benefits id 50544	#9 title Fire with Fire id 139567	#10 title Friends with Kids id 80038					
Ver. v1.1 Hits: 84	title The Basketball Diaries id 10474	title Alien id 348	title Aliens id 679	title Alien³ id 8077	title Cowboys & Aliens id 49849	title Aliens in the Attic id 20856	title Alien: Resurrection id 8078	title Monsters vs Aliens id 15512	title Aliens vs Predator: Requiem id 440	title AVP: Alien vs. Predator id 395					



```
[INFO] --- rre-maven-elasticsearch-plugin:6.3.2:evaluate (search-quality-evaluation) @ rre-fi
[INFO ] RRE: New evaluation session is starting...
[INFO ] RRE: Search Platform in use: Elasticsearch
[INFO ] RRE: Starting Elasticsearch...
[INFO ] RRE: Elasticsearch Search Platform successfully started.
[INFO ] RRE: found 1 ratings sets.
[INFO ] RRE: Ratings Set processing starts
[INFO ]
[INFO ] *****
[INFO ] RRE: Index name => tbmd
[INFO ] RRE: ID Field name => _id
[INFO ] RRE: Test Collection => /Users/agazzarini/workspaces/fosdem/rre-relevant-search-walk
[INFO ] RRE: Loading the Test Collection into Elasticsearch, configuration version v1.0
[INFO ] RRE: Loading the Test Collection into Elasticsearch, configuration version v1.1
[INFO ] RRE: Loading the Test Collection into Elasticsearch, configuration version v1.2
[INFO ] RRE: Elasticsearch has been correctly loaded.
[INFO ] RRE: target versions are v1.0,v1.1,v1.2
[INFO ] TOPIC: Space Jam
[INFO ]   QUERY GROUP: See "Relevant Search", Chapter 3
[INFO ]     QUERY: basketball with cartoon aliens
[INFO ] TOPIC: Star Trek
[INFO ]   QUERY GROUP: Starring
[INFO ]     QUERY: patrick stewart
[INFO ]   QUERY GROUP: Starring and Title
[INFO ]     QUERY: star trek patrick stewart
[INFO ]   QUERY GROUP: Names and Title
[INFO ]     QUERY: star trek patrick stewart william shatner
```

Iteration #3 (v1.2)

- replaces title¹⁰ with title^{0.1}
- introduces shingles for names search
- “Space Jam” search is green!
- Other queries are better as well

RRE DEMO: Iteration #1 (2/2)



FOSDEM 2019

Ranking Evaluation Report - created on Sunday, January 27, 2019

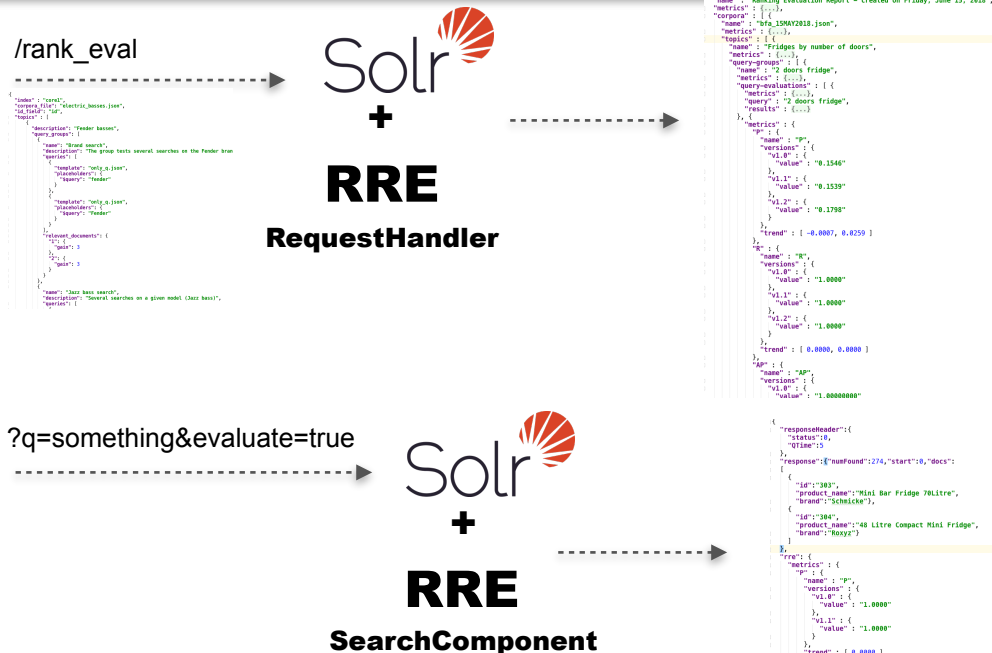
Corpus	Topic	Query Group	Query	Metrics																
				P@1				P@10				AP				NDCG@10				
tmbd.bulk				v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	
				0.5000	0.0000	0.5000	-0.5000	0.3750	0.3750	0.3251	0.0000	-0.0499	0.2853	0.2448	0.4384	-0.0405	0.7126	0.6725	0.8576	-0.0401
Space Jam				v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	
				0.0000	0.0000	1.0000	0.0000	1.0000	0.0000	0.0000	0.1001	0.0000	0.1001	0.0000	0.0000	1.0000	0.0000	0.0000	1.0000	0.0000
		See "Relevant Search", Chapter 3		v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	
				0.0000	0.0000	1.0000	0.0000	0.0000	0.0000	1.0001	0.0000	0.0000	0.1001	0.0000	0.1001	0.0000	0.0000	1.0000	0.0000	
			basketball with cartoon aliens	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	v1.0	v1.1	v1.2	Δ	
			▼	0.0000	0.0000	1.0000	0.0000	0.0000	0.0000	1.0001	0.0000	0.0000	0.1001	0.0000	0.1001	0.0000	0.0000	1.0000	0.0000	
	#1	#2	#3	#4	#5	#6	#7	#8	#9	#10										
Ver. v1.0 Hits: 1197	title Aliens	title The Basketball Diaries	title Cowboys & Aliens	title Monsters vs Aliens	title Aliens vs Predator: Requiem	title Aliens in the Attic	title Dances with Wolves	title Friends with Benefits	title Fire with Fire	title Friends with Kids										
	id 679	id 10474	id 49849	id 15512	id 440	id 20856	id 581	id 50544	id 139567	id 80038										
Ver. v1.1 Hits: 84	title The Basketball Diaries	title Alien	title Aliens	title Alien³	title Cowboys & Aliens	title Aliens in the Attic	title Alien: Resurrection	title Monsters vs Aliens	title Aliens vs Predator: Requiem	title AVP: Alien vs. Predator										
	id 10474	id 348	id 679	id 8077	id 49849	id 20856	id 8078	id 15512	id 440	id 395										
Ver. v1.2 Hits: 84	id 2300	title Grown Ups	title Speed Racer	title Semi-Pro	title The Flintstones	title Coach Carter	title White Men Can't Jump	title Meet Dave	title Aliens vs Predator: Requiem	title Bedazzled										
	title Space Jam	id 38365	id 7459	id 13260	id 888	id 7214	id 10158	id 11260	id 440	id 1636										



- **Search Quality Evaluation**
- **Rated Ranking Evaluator**
- **Demo**

✓ **Future Works / Idea**

- **Q&A**



Rank Eval API

The **RRE** core can be used for implementing a **RequestHandler** which will be able to expose a **Ranking Evaluation** endpoint.

That would result in the same functionality introduced in **Elasticsearch 6.2** [1] with some differences.

- rich tree data model
- metrics framework

Note that in this case it doesn't make so much sense to provide comparisons between versions.

As part of the same module there could be a **SearchComponent** for evaluating a single query interaction.

[1] <https://www.elastic.co/guide/en/elasticsearch/reference/6.2/search-rank-eval.html>



Jenkins Plugin

RRE Maven plugin already produces the **evaluation data** in a **machine-readable format** (JSON) which can be consumed by another component.

The **Maven RRE Report plugin** or the **RRE Server** are just two examples of such consumers.

RRE can be **already executed** in a **Jenkins CI build cycle** (using the **Maven plugin**).

By means of a dedicated **Jenkins plugin**, the evaluation data could be graphically displayed in the Jenkins dashboard. It could be even used for blocking builds which produce bad evaluation results.

[Back to Dashboard](#)[Status](#)[Changes](#)[Workspace](#)[Build Now](#)[Delete Project](#)[Configure](#)[Performance Trend](#)[Build History \(trend\)](#)

#10 Mar 22, 2010 10:36:48 AM

#9 Mar 22, 2010 9:59:28 AM

#8 Mar 22, 2010 9:46:45 AM

#7 Mar 22, 2010 9:38:15 AM

#6 Mar 9, 2010 1:22:57 PM

#5 Mar 9, 2010 12:09:36 PM

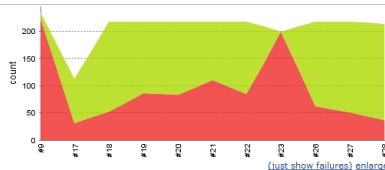
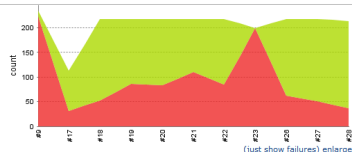
#3 Mar 9, 2010 11:06:32 AM

#2 Mar 9, 2010 10:47:59 AM

#1 Mar 9, 2010 10:38:19 AM

[for all](#) [for failures](#)Last Report
Filter trend data

Topic	Query Group	Query	Metric	NDCG@10										DELTA	
				0.8	0.9	1.1	1.2	1.3	1.4	1.5					
SIM issues				0.4372	0.36	0.9244	0.9244	0.9244	0.9244	0.8844	-0.0772	0.5644	0	0	-0.04
				0	1	1	1	1	1	1	1	0	0	0	0
	The group tests several variants of the topic above.	Sim not readable		0	1	1	1	1	1	1	1	0	0	0	0
		SIM unreadable		0	1	1	1	1	1	1	1	1	0	0	0
Long distance calls				0	1	1	1	1	1	1	1	0	0	0	0
	The group tests several variants of the topic above.	Long distance calls		0	1	1	1	1	1	1	1	0	0	0	0
		Long distance models		0	1	1	1	1	1	1	1	1	0	0	0
	Samsung Galaxy S7 manual				0.5134	0	1	1	1	1	1	-0.5134	1	0	0
The group tests several variants of the topic above.		Tutorial for S7		0.5134	0	1	1	1	1	1	-0.5134	1	0	0	0
		S7 manuals		0.68	0	1	1	1	1	1	1	-0.68	1	0	0
		S7 Manuals		0.43	0	1	1	1	1	1	1	-0.43	1	0	0
Unlock S7				0.43	0	1	1	1	1	1	-0.43	1	0	0	0
				1	0	1	1	1	1	1	-1	1	0	0	0
	The group tests several variants of the topic above.	Unlock S7		1	0	1	1	1	1	1	-1	1	0	0	0
		Unlock Galaxy S7		1	0	1	1	1	1	1	1	-1	1	0	0
	Unlock Samsung S7		1	0	1	1	1	1	1	1	-1	1	0	0	0
	Unlock Galaxy Samsung S7		1	0	1	1	1	1	1	1	-1	1	0	0	0
Unlock (Specific for S7 Edge)				0.75	0	1.25	1.25	1.25	1.25	1.25	-0.75	1.25	0	0	0
	The group tests several variants of the topic above.	Unlock S7 edge		0.75	0	1.25	1.25	1.25	1.25	1.25	-0.75	1.25	0	0	0
		s7 edge unlock		0.75	0	1.25	1.25	1.25	1.25	1.25	-0.75	1.25	0	0	0
				0	1	1	1	1	1	1	1	-1	1	0	-1
Screen Frozen				1	0	1	1	1	1	1	0	-1	1	0	-1
	The group tests several variants of the topic above.	screen frozen can't power off		1	0	1	1	1	1	1	0	-1	1	0	-1
		screen frozen		1	0	1	1	1	1	1	0	-1	1	0	-1
				0	1	1	1	1	1	1	1	0	0	0	0
Cancel Auto-refill				0	1	1	1	1	1	1	1	0	0	0	0
	The group tests several variants of the topic above.	Cancel auto refill		0	1	1	1	1	1	1	1	0	0	0	0
		cancel autorefill		0	1	1	1	1	1	1	1	0	0	0	0
		cancel auto-refill		0	1	1	1	1	1	1	1	0	0	0	0
	how to cancel auto-refill		0	1	1	1	1	1	1	1	1	0	0	0	0
	how to cancel autorefill		0	1	1	1	1	1	1	1	1	0	0	0	0





The main input for RRE is the Ratings file, in JSON format.

Writing a comprehensive JSON to detail the ratings sets for your Search ecosystem can be expensive!



Judgement Collector UI

1. Explicit feedback from users judgements
2. An intuitive UI allow judges to run queries, see documents and rate them
3. Relevance label is explicitly assigned by domain experts

Users Interactions Logger

1. Implicit feedback from users interactions (Clicks, Sales ...)
2. Log to disk / internal Solr instance for analytics
3. Estimate $\langle q, d \rangle$ relevance label based on Click Through Rate, Sales Rate



Once you collected the ratings, could we use them to actively improve the quality metrics ?

Learning To Rank

“Learning to rank is the application of **machine learning**, typically **supervised**, **semi-supervised** or **reinforcement learning**, in the construction of **ranking models** for **information retrieval** systems.” Wikipedia

Users Interactions
Logger

Judgement Collector
UI



Interactions



Trained Model

$$y = f(x)$$

Training



```
1 101:1.2 102:0.03
0 1:2.1 10001:300 10002:400
0 0:1.3 1:0.3
1 0:0.01 1:0.3
0 0:0.2 1:0.3
```




Creating a Learning To Rank Training Set from the collected interactions is not going to be trivial. It normally requires ad hoc data manipulation depending on the use case...
... but some steps could be automated and make available for a generic configurable approach

Configuration

- Null feature sanitisation
- Query Id calculation
- Query document feature generation
- Single/Multi valued categorical feature encoding



1. Ad Hoc category, Artificial values, keep NaN
-> depends of Training Library to use
2. Optional Query Level features to be hashed as QueryId
3. Intersect related query and document level categorical features to generate Ordinal query-document features
4. Label Encoding ? One Hot Encoding? Binary Encoding? [1]
 Dummy Variable Trap



What about the relevance label for each training vector ?
Can we estimate it from the interactions collected ?

- Interaction Type Counts
- Click Through Rate/Sales Through Rate calculation
- Relevance label normalisation



Configuration

1. Impressions? Clicks? Bookmarks? Add To Charts? Sales?
2. Define the objective: Clicks/Impressions ? Sales/Impressions?
3. Relevance Label : 0...4



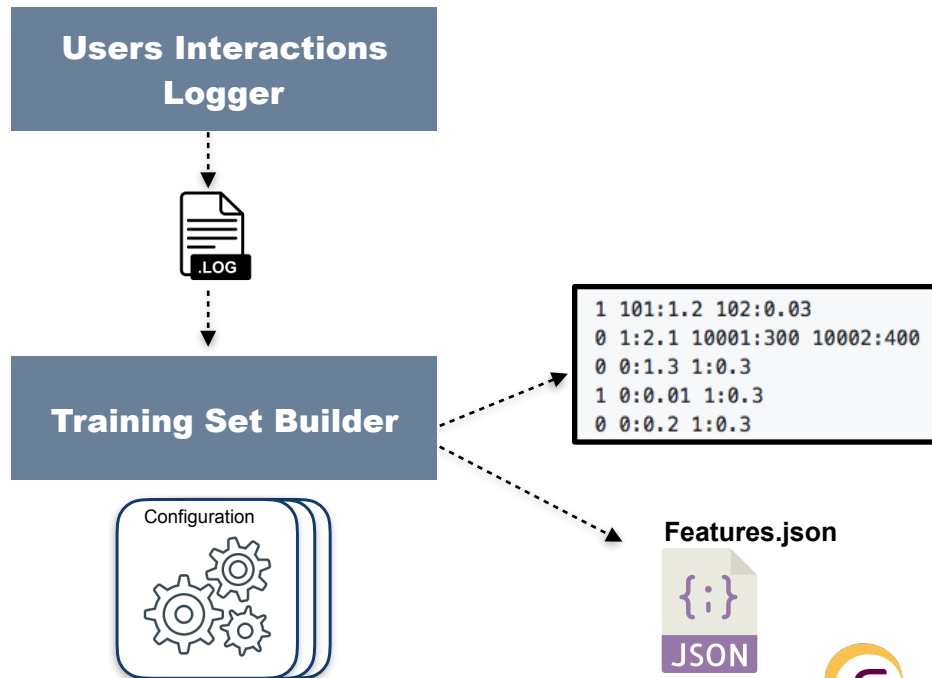
Can the features.json configuration generation be automated?

Learning To Rank - Solr features.json

The features.json is a configuration file necessary for Solr Learning To Rank extension to work.

It is a configuration file that describes how the features that were used at training time for the model can be extracted at query time.

This file is coupled both with the training set features and the query time features.





- **Search Quality Evaluation**
- **Rated Ranking Evaluator**
- **Future Works**
- **Demo**

✓ **Q&A**



FOSDEM 2019



RRR

Rated Ranking Evaluator

An open-source approach for Search Quality Evaluation

Alessandro Benedetti, Software Engineer

Andrea Gazzarini, Software Engineer

3rd February 2019

Thank you!



Sease