

Sease

Search Quality Evaluation

Tools and Techniques

Alessandro Benedetti, Software Engineer

Andrea Gazzarini, Software Engineer

2nd October 2018

Who we are

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

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





Who we are

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

Andrea Gazzarini, "Gazza"

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







🕭 Sun

을 Java



≥Sui

▶Sun

≥Sun



Sease

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

2026

Search Services

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



Agenda

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

✓ Search Quality Evaluation

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



Search Quality Evaluation



Context Overview

Search engineering is the production of quality search systems.

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

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.



Search Quality Evaluation: Correctness



Correctness

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

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

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.

Search Quality Evaluation / Measures

We are mainly focused here

Evaluation Measures

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

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**.

Agenda

Search Quality Evaluation

✓ Rated Ranking Evaluator (RRE)

- What is it?
- How does it work?
- Evaluation Process Input & Output
- Challenges
- > Future Works
- > Q&A

RRE: What is it?

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

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

Difference of the second secon	* +- 🕱-				
SeaseLtd / rated-ranking-evaluator	3 ★Star 16 ¥fork 0				
O Code ⊙ Issues (8) Pull requests @ Projects @ Wild \	ĝs.				
Home Andrea Gazzarini edited this page on Jul 16 - 9 revisions	Edit New Page				
BRRE	+ Pages 🗃				
Rated Ranking Evaluator	2. Quick Start 3. Project Structure 4. Evaluation Measures 5. How does it work?				
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.	5.1 Domain Model 5.2 What we need to provide 5.3 Where we need to provide 5.4 The Evaluation Process				
It is something which helps a Search Engineer in his daily job. Are you a Search Engineer? Are you	5.5 The Evaluation Output 6. RRE Server				

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" loval providing human reactable outputs that could be act as deliverables for

RRE: What is it?

- A set of search quality evaluation tools ٠
- A search quality evaluation framework ٠
- Multi (search) platform ٠
- Written in Java •
- It can be used also in **non-Java projects** ٠
- Licensed under Apache 2.0 ٠
- Open to contributions ٠
- Extremely dynamic! ٠

RRE: At a glance

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

RRE: Ecosystem

HAYSTACK THE SEARCH RELEVANCE CONFERENCES

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.

RRE: Available metrics

- 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
- F-Measure >>>> Compound Metric

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

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 (1/2)

RRE Domain Model

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

HAYSTACK

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 (2/2)

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 possibile 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)

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

RRE: Evaluation process overview (2/2)

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

RRE: Corpora

HAYSTACK THE SEARCH RELEVANCE CONFERENCES

Solr®

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

elasticsearch

{"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**.

RRE: Configuration Sets

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

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)

RRE / Query templates

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**.

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

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)

RRE: Ratings

Ratings

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

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

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.

RRE: Evaluation Output

Evaluation output

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

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

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

RRE: Workbook

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

			c	0		#	6	-		- 2		L MNO	P	0	
5	Texale	Overy Crown	Owers	Matheir											
÷	- Aller	den la receb	sport y					M	00.0010						
ŀ							1.0		1.4			10115			
÷				0.000	1.10	100004	1.0	10000	10000	1.0		01100 0 0 0 0		0.000	
è	Table secures										-		-		
è		The second basis and	and contacts of the topic should		- 2	- 2	- 2	- 2	- 3	-					
2			fam not manafelia	_	- 2	- 2	- 2		- 2	_					
ï			Tild completing			- 2	- 2								
ì	Long distance calls				- 1	- 1	- 1	- 1	- 1	-				0	
k		The group leads any	mode sign with a strainer later.		- 1		- 2			- 1					
Ē			Long-distance calls											0	
ā	1		Long distance models		- 1	- 1		- 1	- 1		- 1			0	
i	Sansurg Galaxy S7 nanual			0.5534	- 0	1			- 1	-	-0.5136			0.5354	
í		The group hosts any	works sign with a strain when	0.55(94			- 2	- 6		- 1	-0.5124			0.03884	
ŝ	1		Tutorial for ST	0.68	0	1		1	1	1	-0.68	1000		0.4	
i			17 menuals	0.43	0	1	1	1	1	1	4.43			0.8	
í			57 Manuals	0.43	- 0				1		4.43			0.6	
i	luminok SP			8		1	1		1		- 4			0.5	
į		The group tests are	veral variants of the topic above.	1	0	1	1	1	1	1	- 4	1		0.1	
	1		United B7	1	0	1	1	1	1	1	- 4			0.1	
			United Galaxy S7		- 0	1			1					0.5	
į	1		United Samurg 57	8	- 0	1	1	1	1		- 4			0.5	
1	1		United Galaxy Samaung 57	1	0	1	1	1	1	1	- 4	1 4 4 4		0.1	
į	United (Specific for \$7 Edge)			0.75	0	1.35	1.25	1.25	1.35	1.25	-4.75	1.25 0 0 0		0.1	
į		The group tests sev	retal variants of the topic above.	0.75	- 0	1.25	1.25	1.25	1.25	1.25	4.75	1.25 4 4 4		0.5	
į	1		United S7 edge	0.75	- 0	1.05	1.05	1.85	1.05	1.05	4.75	1.25		0.5	
į			s7 edge-unicox	0.75	0	1.29	1.25	1.25	1.25	1.25	-4.75	1.25 0 0 0		0.1	
į	Roman Property			1	0	1	1	1	1		- 4		- 4	1	
	1	The group tests sev	retal variants of the topic above.		- 0	1			1					1	
			scene from cart power of			1			1		- 4		- 4	1	
j			scelet frozen	1	0	1	1	1	1		- 4	1000	- 4	1	
	Canoel Auto nelli				1	1	1	1	1					0	
į		The group tests sev	retral variants of the topic above.		1				1					- 0	
Í			Carear auto ratio		1				1						
ĺ			canod autorette		1	1	1	1	1	1	1			0	
			cancel auto refil		1	1		1	1	1				0	
ł			how to cancel auto-refit		1										
ļ			how to cancel autoratid		1				1						
į	Text nessaging with VV/O 5			0.7467	0	1	1	1	1	1	-0.7467	1000		0.8	
l		The group tests say	renal variants of the topic above.	0.7467	0	- 1			1	- 1	4.7407			0.3	
			Text messaging with VIVD-5	0.90					1		-0.93			0.4	
l			held message VIVO 5						1		- 4			0.4	
ļ			she work	0.30	0	1		1	1	1	4.8			0.1	
ł	Text messaging (send) VIVO			0.75	0	1		1	1		4.75			0.1	
į		The group tests sev	veral variants of the topic above.	0.75					1		4.75			0.5	

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**

RRE: RRE Server (1/2)

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

		-						
ing Colorador	Report -		e Medweley Are 21.	378				
Orque	744	8.41) 0.140	Barry			8471		Matrice Andrews
54072978,344					24 20 20 4	10.0 00.0 10.0 0	108 101 108 8	
					100 100 1000 000 000	44.8 44.8 44.8 8 1000 1000 1000 1000	418 45 45 8 1000 1000 1000 1000	
		3 dere Holge		48 48 48 8	*** *** ***	*** *** *	We do we have	*** *** ***
			a mark frame			*****	*** *** **	***
			k mer fridge k	48 48 48 8	418 418 418 8	418 415 415 8 (1000 1000 1000 1000	with with with the Content Content Content Content	418 413 413 8 10000 10000 1000 1000
			Star 1-star	VA AA VA A	104 00 004 4	AND AN AND A	VA AN AN A	244 253 254 A
			ittern frohe	VA VA VA A	VIA VIA VIA A	44 45 44 4	VII 43 VII 4	44 45 45 4
	104			44 45 45 4		era era era a	100 00 00 000 000	and and and a
		111		418 41 418 8	eta eta eta a Canto Datos Loto Loto I	elle ell elle a	ela el ela a TREP ESEP DEES LOOM 11	eta eti eta a Eleti ituti ituti anni itu
			Ber Folge sont offer		eta eta eta a Canto Canto Vitto anno 1	48 48 48 8 1107 1107 1200 1000 1000	elà elà elà à Esta d'aca della <mark>dese</mark> la	#48 #61 #18 &
	Coloured Anna Lingue				204 00 000 4	104 00 00 4	214 21 23 4	
		-		VIA 417 418 8	V18 113 V18 8	44 45 43 A	10.0 10.0 10.0 0	418 41 418 B

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: RRE Server (2/2)

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

The evaluation data, at query / version level, collects the top n search results.

In the web console, under each query, there's a little arrow which allows to open / hide the section which contains those results.

In this way you can get immediately the meaning of each metric and its values between different versions.

In the example above, you can immediately see why there's a loss of precision (first metric) between v1.0, v1.1, which got fixed in v1.2

RRE: Iterative development & tuning

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

We are thinking about how to fill a third monitor

RRE: Challenges

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

"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?

Java is not in my stack

Agenda

- Search Quality Evaluation
- > Rated Ranking Evaluator

✓ Future Works / Idea

> Q&A

Future Works: Solr Rank Eval API

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

Future Works: Jenkins Plugin

Back to Destroyed Q. Stetus Last Report Filter loand data Changes Query Group Query Matrix Tertapace Teople NDCG@10 Dubl.Nor 4.4 1.3 1.5 1.6 L4 1.664 8:5644 · S Delete Project DW MARK The price leafs several variants of the bass above. Configure Que and employed line presidents E Performance Trend Long distance cafe he prove texts several variants of the tops above Build History (band) Long distance calls Long-Bistance models Nor 22, 2010 10:36:48 AM Sahaung Dates 37 Awarust 0.00.00 ler 22, 2010 9-59-28 AM The group limits serviced variants of the topic above 0.62.76 6.1124 Leonial Nor 57 0.66 4.68 Ner 22, 2010 9:46:45 AM 0.40 5P nanuals for 22, 2010 5-38-15 AM 17 Manuals 0.43 United 17 Ser 9, 2010 1-22-57 PH group tests service' variants of the task above. Nor 9, 2010 12:09:36 PW second and United Salary 87 Har 9, 2010 11-06-12 AM United Gampung C Q #2 Her 9, 2010 10:47/39 AM unitoti Galary Samaing S² princip (Taperilly for \$17 billion 0.75 1.75 1.75 1.00 1.75 #1 Nor 5, 2010 10:38:13 AM he proof book as mini variante of the logic along 1.04 1.95 1.04 1.04 1.00 1.84 4.76 the all the failures United STratigs 1.25 1.25 1.24 1,255 1.24 of white unlock 1.05 Annual Property group into several variants of the topic above sensor from our's many of screet frank Canal Auto-telli The group lexits serving variants of the topic above Carrow and well Carried Bulleville second pulse will how to carved mate call Now to cancel autorable

Jenkins Plugin

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

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

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.

Future Works: Building the input

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 <q,d> relevance label based on Click Through Rate, Sales Rate

HAYSTACK

THE SEARCH RELEVANCE CONFERENCE

Future Works: Learning To Rank

HAYSTACK THE SEARCH RELEVANCE CONFERENCE!

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

Future Works: Training Set Building

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

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

Configuration

- Ad Hoc category, Artificial values, keep NaN

 > depends of Training Library to use
- 2. Optional Query Level features to be hashed as QueryId

HAYSTACK

- 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

[1] https://www.datacamp.com/community/tutorials/categorical-data

Future Works: Training Set Building

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

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

Future Works: Learning To Rank Solr Configs

HAYSTACK THE SEARCH RELEVANCE CONFERENCE

Can the features.json configuration generation be automated?

- > Search Quality Evaluation
- Rated Ranking Evaluator
- > Future Works

√ Q&A

