

Learning to Boost

Logistic Regression to Optimize Elasticsearch Boosts

Nina Xu and Jenna Bellassai Haystack Conference 2021

What is Guru?

Guru is a company wiki that works in your workflow.

With Guru, you can asynchronously:

- Share critical product information
- Onboard employees autonomously
- Streamline internal communications



Helping the most innovative companies in the world work smarter—from anywhere.

ZOOM # slack shopify noom

Create and share trusted content



Who We Are

5

Search & Discovery Team





Tracy Hall FE Eng



Sean Fleming BE Eng









Nabin Mulepati ML Eng



Bernie Gray Data Science



Yev Meyer Data Science



Jake Sauer Design



Adam Gruber FE Eng



Bennett Ferris Data Science





Nina Xu Data Science

Search at Guru

- Customers search for information contained in "Cards" in their own instance
- B2B use case: information is particular to each customer
- Elasticsearch



Information landscape

- Thousands of companies
 - From a variety of industries
 - Create and maintain their own documents
 - Company-specific jargons
 - Customers are the subject matter experts, not us
- Large companies
 - 10k+ documents
 - 5k+ of queries per day
- Mid-sized companies
 - ~100+ documents
 - A few queries per day



Outline

- 1. Problem Statement
- 2. Approach
- 3. Results
- 4. Special Considerations

Problem Statement

Elasticsearch query (simplified)





Elasticsearch query (realistic)

Many field boosts to tune

		ļ
3-	"bool": {	
4 -	"must": L	
5 -	1	
6 -	"bool": 1	
7.	"must":	
8 -	1	
9 +	"bool": {	
10 -	"should": [
11 -	{	
12 -	"multi_match": {	
13	"query": "rocky",	
14 -	"fields": [
15	"title^10",	
16	"title.text_std",	
17	"title.text_std_gr_idx_syn",	
18	"title.text_std_syn",	
19	"title.text_std_idioms",	
20	"cast^2",	
21	"cast.text_std",	
22	"cast.text_all",	
23	"cast.text_people",	
24	"cast.text_people_max",	
25	"cast.text_people2",	
26	"cast.text_people_notf",	
27	"cast.text_std_gr_idx_syn",	

Review: how to choose Elasticsearch boost values?

Manual tuning



Need Better Search? Take Control with Quepid

Built by search experts at OpenSource Connections, Quepid puts your team in control of your Soir & Elasticsearch search results. Using test-driven techniques, content experts and marketing teams can efficiently define and track search correctness. With this foundation, you'll never slide backwards and can fully leverage the power of open source search.

Quepid is available as a free, hosted service - try it now - and is built on open source software allowing you to download and build it into your own infrastructure.

auepid.com/

Read About Quepid's Simple Secret »

- Advantages:
 - Easy to use
 - Guided by search metrics
- Limitations:
 - Requires explicit relevance judgements
 - Requires some domain expertise
 - Hard to extend to multi-tenant use case

Grid search

9 🗗 in



- • If the overall quality *increases* and meets your quality SLAs, you're done, otherwise repeat until you reach the desired quality.
- If the overall quality decreases, revert the change, make a note of it for future reference and try something else.

elastic.co/blog/test-driven-relevance-tuning-of-elas ticsearch-using-the-ranking-evaluation-api

- Advantages
 - More thorough than "try it and see"
- Limitations
 - Permutation explosion

Genetic algorithms

HAYSTACK

me Schedule Tickets Training Travel Past Conferences Past Tall

Evolving Relevance

im Allison • Location: Theater 4 • Back to Haystack 2020



This talk builds on work by Simon Hughes and others to apply genetic algorithms (GA) and random search for finding optimal parameters for relevance ranking. While manual tuning can be useful, the parameter space is too vast to be confident that one has found optimal parameters without overfitting. We'll present Quaerite (https://github.com/tballison/quaerite), an open source toolkit that allows users to specify experiment parameters and then run a random search and/or a GA to identify the best settings given ground truth.We'll offer an overview of mapping parameter space to a GA problem in both Solr and Elasticsearch, the importance of the baked-in n-fold cross-validation, and the surprises and successes found with deployed search systems.

haystackconf.com/us2020/evolving-relevance/

<u>github.com/tballison/quaerite</u>

Advantages

- Data-driven
- Can test many parameters besides just boosts
- No linearity constraint (more on this later)
- Limitations
 - May not scale well with complexity

Learning to Rank



4/elasticsearch-learning-to-rank/

github.com/o19s/elasticsearch-learning-to-rank

- Advantages:
 - Data-driven
- Limitations:
 - High data need
 - High barrier to entry
 - Usually done at reranking step due to high computation demand

Approach

Learning to Boost - How It Works

Implementation



Learning to Boost (LTB) is a logistic regression model that uses relevance judgements to determine the optimal Elasticsearch boost values for an Elasticsearch query.

Why should you care?

- Data-driven
- Easy to train
- Easy to productionize (you'll see)
- Automated for future iterations



How it works - data requirement



જિ

Elasticsearch scoring (commonly)



Optimal boost





Elasticsearch scoring (commonly):

Logistic regression:

$$\ell = \log_b rac{p}{1-p} = eta_0 + rac{eta_1 x_1}{p} + rac{eta_2 x_2}{p}$$

where p is the probability that a binary label is 1

Implementation





Implementation

Collect Data

Replay past searches with all boosts set to 1. Record explanations for each result.

How do we replay past searches?

- We use a homegrown framework ("offline search trials") to replay millions of past searches using a specified algorithm in an environment isolated from production.
- We compare the results of those replayed searches to the results we saw in prod.
- For Learning to Boost, an initial search trial using an Elasticsearch query with boosts set to 1 provides us with training data (search explanations and scores).

Collect Data





Collect Data



Collect Data

Implementation

Collect Data

Featurize

Replay past searches with all boosts set to 1. Record explanations for each result. Transform search scores and explanations from trial into features. Add labels according to user behavior.

31

Featurize

{

value: 16.460304, description: "sum of:",

- details: [- {

```
search
            doc_id
                       title
                                   overvi
                                               label
_id
                                   ew
. . .
            . . .
                        . . .
                                   . . .
                                               . . .
172
                       8.838
                                   7.622
            2
                                               1
```

. . .

. . .

. . .

. . .

. . .

```
value: 8.837944,
     description: "weight(title:rocki in 1335) [PerFieldSimilarity], result of:",
    - details: [
      - {
           value: 8.837944,
           description: "score(freq=1.0), computed as boost * idf * tf from:",
          + details: [...]
                                                                                     {
        3
     1
                                                                                           "total_score": 16.460304,
  },
                                                                                           "title": 8.837944,
 - {
     value: 7.6223593,
                                                                                           "overview": 7.6223593
     description: "weight(overview:rocki in 1335) [PerFieldSimilarity], result of:",
    - details: [
                                                                                     }
       - {
           value: 7.6223593,
           description: "score(freq=2.0), computed as boost * idf * tf from:",
          + details: [...]
        }
     1
  3
1
```

Implementation

Replay past searches with all boosts set to 1. Record explanations for each result. Transform search scores and explanations from trial into features. Add labels according to user behavior. Train logistic regression model.

The pairwise approach

G

Why?

Elasticsearch scores are not comparable across queries

The pairwise approach

Pointwise LTR:

Pairwise LTR:

Pairwise data

How?

- Take difference of the feature values for each pair of docs from the same query
- Create new labels based on comparison of relevance

Logistic regression on pairwise data

Let *i*, *j* denote two documents from the search result list.

If $l_i > l_j$, doc *i* is more relevant than doc *j*

Fit logistic regression model

Use your favorite package to make it happen!

- pyspark.ml
- scikit-learn

Make sure to restrict coefficients to non-negative values

Implementation

Replay past searches with all boosts set to 1. Record explanations for each result. Transform search scores and explanations from trial into features. Add labels according to user behavior. Train logistic regression model.

Run and evaluate a search trial that replays past searches using learned model coefficients as boost values.

Implementation

Deploy

Just change the boost values.

What if we need to add a new field?

What if we need to add a new field?

Results

Model	AUC	MAP@5
Baseline (hand-tuned)	0.897	0.332
LTB	0.911 (+1.6%)	0.336 (+1.2%)

Special Considerations

Some prerequisites (for us)

- Event tracking enabled us to obtain *implicit* judgements
- Our home-grown offline search trial framework allowed us to
 - Obtain training data
 - Replay searches with new boost values
 - Calculate search metrics

Linearity constraint

The regression approach is **useful for**

- Queries where field scores are summed together
- Identifying top-level boost values if granular-level field scores are not additive

The regression approach is **not useful for**

- the script_score portion of the query*
- multi-match queries with best_fields type
- dis_max queries with tie_breaker

Labeling considerations

Implicit vs. explicit judgment

- Both work
- Implicit judgment allows for bigger sample size and is suitable for the enterprise use case
- Explicit judgment is less prone to position bias

Binary vs. graded relevance

- Use logistic regression for binary labels
- Use ordinal regression for graded labels eg. 1, 2, 3, 4

Loss function vs. search metric

- The binary logistic loss (cross-entropy loss) does not take into account position of documents
- Ranking metrics in search usually consider position
- In practice, the logistic regression metrics & ranking metrics generally agree until AUC gets very close to 1

LTB is not intended to replace LTR 😊

Presentation resources

• See it in GURU

We are hiring!

getguru.com/careers

...

- Sr. Search Engineer
- Sr. Machine Learning Engineer
- Sr. Full Stack Engineer
- Sr. Backend Engineer
 - Sr. Frontend Engineer

Don't take yourself too seriously Seek balance. Assume good intent.

Seek and share knowledge

Be transparent. Seek diverse views for better outcomes.

Based in Philadelphia, San Francisco, or remote!

Thank you!