# Elastic Search: Similarity & Scoring[1–4]

How Retrieved Documents are Ranked

Weimao Ke

Drexel University

# Table of contents

# Mapping and Built-in Similarity

## Mapping Similarity

With Mapping:

- You can configure a scoring or similarity function for each field.
- Basic similarity functions include: *BM25* (default), *classic* (TF\*IDF), and *boolean*
- Custom similarities can be configured with related parameters.

This changes the way in which retrieved documents are ranked (sorted).

It is the most important function of a search engine.

## Mapping Similarity

Example:

```
1    PUT my_index
2    {
3      "mappings": {
4        "properties": {
5          "field1": {
6            "type": "text"
7          },
8          "field2": {
9            "type": "text",
10           "similarity": "boolean"
11         }
12       }
13     }
14   }
```

- By default, *field1* uses *BM25* similarity.
- *Field2* is configured to use the *boolean* similarity.

## Mapping Similarity

You can change the default similarity for a new index:

```
PUT /index
{
  "settings": {
    "index": {
      "similarity": {
        "default": {
          "type": "boolean"
        }
      }
    }
  }
}
```

This is done when the index is created.

## Mapping Similarity

You can change the default similarity of an existing index.

You will have to close the index to change it:

```
POST /index/_close

PUT /index/_settings
{
  "index": {
    "similarity": {
      "default": {
        "type": "boolean"
      }
    }
  }
}
```

And reopen it after the change:

```
POST /index/_open
```
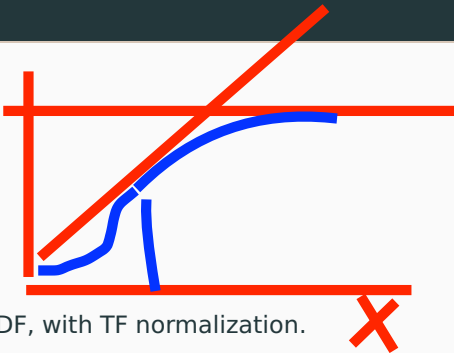
# Custom Similarities

BM25 similarity (default):

- Regarded as a version of TF*IDF, with TF normalization.
- $k_1$: controls TF normalization (saturation), default 1.2;
- $b$: controls document length normalization, default 0.75.

## Similarity Modules

DFR (Divergence from Randomness) similarity:

- *basic_model*: possible values *g*, *if*, *in*, and *ine*;
- *after_effect*: possible values *b* and *l*;
- *normalization*: possible values *no*, *h*1, *h*2, *h*3, and *z*.

All options but the first option need a normalization value.

## Similarity Modules

Other similarity functions include:

- DFI (Divergence from Independence) similarity
- IB (information based model) similarity
- LM Dirichlet similarity
- LM Jelinek Mercer similarity

## Mapping Similarity

One can create a customized version of a built-in similarity function:

```
1   PUT /index
2   {
3     "settings": {
4       "index": {
5         "similarity": {
6           "my_similarity": {
7             "type": "DFR",
8             "basic_model": "g",
9             "after_effect": "l",
10            "normalization": "h2",
11            "normalization.h2.c": "3.0"
12          }
13        }
14      }
15    }
16  }
```

And then use the custom version with mapping:

```
1   PUT /index/_mapping
2   {
3     "properties" : {
4       "title" : { "type" : "text", "similarity" : "
            my_similarity" }
5     }
6   }
```

## Mapping Similarity

ElasticSearch also allows scripted similarity:

```
1   PUT /index
2   {
3     "settings": {
4       "number_of_shards": 1,
5       "similarity": {
6         "custom_tfidf": {
7           "type": "scripted",
8           "script": {
9             "source": "
10              double tf = Math.sqrt(doc.freq);
11              double idf = Math.log((field.docCount+1.0)/(term.
                  docFreq+1.0)) + 1.0;
12              double norm = 1/Math.sqrt(doc.length);
13              return query.boost * tf * idf * norm;"
14          }
15        }
16      }
17    },
```

## Mapping Similarity

Scripted similarity (continued):

```
1    "mappings": {
2      "properties": {
3        "field": {
4          "type": "text",
5          "similarity": "custom_tfidf"
6        }
7      }
8    }
9  }
```

## Mapping Similarity

Note:

- *doc.freq* is NOT document frequency (DF), but a term's frequency in the document (TF).
- *term*.*docFreq* is document frequency (DF).

Rules on scripted similarity:

- Returned scores must be positive;
- All other variables remaining equal, scores must not decrease when term frequencies (TF) increases;
- All other variables remaining equal, scores must not increase when document length increases.

## Mapping Similarity

You may also use *weight_script* for more efficient scoring:

```
PUT /index
{
  "settings": {
    "number_of_shards": 1,
    "similarity": {
      "scripted_tfidf": {
        "type": "scripted",
        "weight_script": {
          "source": "double idf =
            Math.log((field.docCount+1.0)/(term.docFreq+1.0)) + 1.0;
            return query.boost * idf;"
        },
        "script": {
          "source": "
            double tf = Math.sqrt(doc.freq);
            double norm = 1/Math.sqrt(doc.length);
            return weight * tf * norm;"
        }
      }
    }
  },
```

## Mapping Similarity

*weight_script* (continued):

```
1    "mappings": {
2      "properties": {
3        "field": {
4          "type": "text",
5          "similarity": "scripted_tfidf"
6        }
7      }
8    }
9  }
```

# Script Score Query

## Script Score Query

Query with a score script:

- *script_score* provides a custom score as part of a query request.
- Changes the way in which returned documents are scored and ranked for the query.

## Script Score Query

Example script score query:

```
1   GET /_search
2   {
3       "query" : {
4           "script_score" : {
5               "query" : {
6                   "match": { "message": "Best Video Tutorials"
                        }
7               },
8               "script" : {
9                   "source" : "doc['likes'].value / 10 "
10              }
11          }
12      }
13  }
```

Suppose *doc*['*likes*'] is a field with the count of likes.

Here it ranks documents according to the number of likes.

## Script Score Query

Parameters for *script_score*:

- *query*: the query to search and return documents;
- *script*: the script used to compute the score of each returned document;
- *min_score* (optional): the threshold score for a document to be included in the search results (hits).

Note that *_score* variable has the current relevance score for each document and can be used as part of the new scoring function.

## Predefined functions for *script_score*

Saturation:

$saturation(value, k) = value/(k + value)$

```
1   "script" : {
2       "source" : "saturation(doc['likes'].value, 1)"
3   }
```

that is $\frac{likes}{1+likes}$:

- 0 when $likes = 0$;
- 0.5 when $likes = 1$;
- 1 when $likes \rightarrow \infty$

# Predefined functions for *script_score*

Sigmoid:

$$sigmoid(value, k, a) = value^a / (k^a + value^a)$$

```
1   "script" : {
2       "source" : "sigmoid(doc['likes'].value, 2, 1)"
3   }
```

that is $\frac{likes^1}{2^1 + likes^1} = \frac{likes}{2 + likes}$

## Predefined functions for *script_score*

Random score:

*randomScore*($< seed >, < fieldName >$)

```
1   "script" : {
2       "source" : "randomScore(100)"
3   }
```

- This generates some randomness (noise) in the ranking, e.g. for random sampling in the search results.
- $< fieldName >$ parameter, as the source of randomness, is optional. Document ID is used by default to generate the random score.

Decay functions:

- Provides scoring based on some kind of distance for numeric values, dates, and geo locations.

An example for geo fields:

```
1   "script" : {
2       "source" : "decayGeoExp(params.origin, params.scale,
            params.offset, params.decay, doc['location'].value)",
3       "params": {
4           "origin": "40, −70.12",
5           "scale": "200km",
6           "offset": "0km",
7           "decay" : 0.2
8       }
9   }
```

## Predefined functions for *script_score*

Functions for vector fields (experimental):

- Functions to be used on *dense_vector* and *sparse_vector* fields.
- A score is computed for every matched document so try to limit the number of docs with a query parameter.

An index with vector fields:

```
 1   PUT my_index
 2   {
 3     "mappings": {
 4       "properties": {
 5         "my_dense_vector": {
 6           "type": "dense_vector",
 7           "dims": 3
 8         },
 9         "my_sparse_vector" : {
10           "type" : "sparse_vector"
11         },
12         "status" : {
13           "type" : "keyword"
14         }
15       }
16     }
```

# Predefined functions for *script_score*

Example vector data:

```
PUT my_index/_doc/1
{
  "my_dense_vector": [0.5, 10, 6],
  "my_sparse_vector": {"2": 1.5, "15" : 2, "50": −1.1, "4545": 1.1},
  "status" : "published"
}

PUT my_index/_doc/2
{
  "my_dense_vector": [−0.5, 10, 10],
  "my_sparse_vector": {"2": 2.5, "10" : 1.3, "55": −2.3, "113": 1.6},
  "status" : "published"
}
```

## Predefined functions for *script_score*

A script score query using cosine similarity:

```
1   GET my_index/_search
2   {
3     "query": {
4       "script_score": {
5         "query" : {
6           "bool" : {
7             "filter" : {
8               "term" : {
9                 "status" : "published"
10              }
11            }
12          }
13        },
14        "script": {
15          "source": "cosineSimilarity(params.query_vector, doc['
                my_dense_vector']) + 1.0",
16          "params": {
17            "query_vector": [4, 3.4, -0.2]
18          }
19        }
20      }
21    }
```

## Predefined functions for *script_score*

Cosine similarity of sparse vectors:

```
 1   GET my_index/_search
 2   {
 3     "query": {
 4       "script_score": {
 5         "query" : {
 6           "bool" : {
 7             "filter" : {
 8               "term" : {
 9                 "status" : "published"
10               }
11             }
12           }
13         },
14         "script": {
15           "source": "cosineSimilaritySparse(params.query_vector, doc['
                 my_sparse_vector']) + 1.0",
16           "params": {
17             "query_vector": {"2": 0.5, "10" : 111.3, "50": −1.3, "113":
                 14.8, "4545": 156.0}
18           }
19         }
20       }
```

## Predefined functions for *script_score*

Other functions for vector fields:

- *dotProduct* computes the dot product of two vectors;
- *dotProductSparse* is the sparse vector version for dot product calculation;
- *l1norm* or *l1normSparse* calculates the $L^1$ distance, i.e. Manhattan distance;
- *l2norm* or *l2normSparse* ccalculates the $L^2$ distance, i.e. Euclidean distance

# Practical Scoring with Relevance Signals

## Relevance Signals

Additional relevance signals:

- Relevance based on terms is not the only source of evidence.
- Other signals correlated with relevance: quality, reputation/popularity, importance, etc.
- Hyperlinks, PageRanks, # of likes, etc. are signals that can be used in relevance scoring (ranking).
- Can be implemented using *script_score* or *rank_feature* queries.

## Relevance Signals

Scoring with PageRank (using *script_score*):

```
1   GET index/_search
2   {
3       "query" : {
4           "script_score" : {
5               "query" : {
6                   "match": { "body": "University" }
7               },
8               "script" : {
9                   "source" :
10                      "_score * saturation(doc['pagerank'].value, 10)"
11              }
12          }
13      }
14  }
```

that is $bm25\_score \times \frac{pagerank}{10+pagerank}$

Here *pagerank* must be mapped as a *numeric* field.

## Relevance Signals

Scoring with PageRank (using *rank_feature*):

```
1    GET _search
2    {
3        "query" : {
4            "bool" : {
5                "must": {
6                    "match": { "body": "University" }
7                },
8                "should": {
9                    "rank_feature": {
10                       "field": "pagerank",
11                       "saturation": {
12                           "pivot": 10
13                       }
14                   }
15               }
16           }
17       }
18   }
```

Here *pagerank* must be mapped as a *rank_feature* field.

Comparing *script_score* vs. *rank_feature*:

- *script_score* is more flexible and offers more ways to combine text relevance with other signals.
- *rank_feature* is limited but more efficient (faster);
- *rank_feature* fields are indexed signals, based on which the query can skip non-competitive documents and find the top matches much faster.

More on *rank_feature*:

- *rank_feature* is typically used in the *should* (OR) clause of a *bool* query;
- this is to *add* the signal to the other scores.
- supports a few functions: saturation, logarithm, and sigmoid.
- must setup *rank_feature* or *rank_features* fields with the index.

## Relevance Signals

Setup an index for *rank_feature*:

```
1   PUT /test
2   {
3     "mappings": {
4       "properties": {
5         "pagerank": {
6           "type": "rank_feature"
7         },
8         "url_length": {
9           "type": "rank_feature",
10          "positive_score_impact": false
11        },
12        "topics": {
13          "type": "rank_features"
14        }
15      }
16    }
17  }
```

*positive_score_impact* : *false* for URL's negative impact on relevance; *rank_features* for the topics field that will contain a list of topics and their relevance.

## Relevance Signals

Index documents with *rank_feature(s)*:

```
 1   PUT /test/_doc/1?refresh
 2   {
 3     "url": "http://en.wikipedia.org/wiki/2016_Summer_Olympics",
 4     "content": "Rio 2016",
 5     "pagerank": 75.3,
 6     "url_length": 42,
 7     "topics": {
 8       "sports": 50,
 9       "brazil": 30
10     }
11   }
12
13   PUT /test/_doc/2?refresh
14   {
15     "url": "http://en.wikipedia.org/wiki/2016_Brazilian_Grand_Prix",
16     "content": "Formula One motor race held on 13 November 2016",
17     "pagerank": 38.9,
18     "url_length": 47,
19     "topics": {
20       "sports": 35,
21       "formula one": 65,
22       "brazil": 20
23     }
24   }
```

## Relevance Signals

Query with *rank_feature*(*s*):

```
1   GET /test/_search
2   {
3     "query": {
4       "bool": {
5         "must": [
6           {
7             "match": {
8               "content": "2016"
9             }
10          }
11        ],
12        "should": [
13          {
14            "rank_feature": {
15              "field": "pagerank"
16            }
17          },
```

## Relevance Signals

Query with *rank_feature*(*s*) (continued):

```
 1              {
 2                "rank_feature": {
 3                  "field": "url_length",
 4                  "boost": 0.1
 5                }
 6              },
 7              {
 8                "rank_feature": {
 9                  "field": "topics.sports",
10                  "boost": 0.4
11                }
12              }
13            ]
14          }
15        }
16    }
```

*boost* (relative to 1.0) is to boost – increase ($> 1$) or decrease ($< 1$)
– a relevance score.

## Relevance Signals

Saturation – example $\frac{pagerank}{pagerank+8}$ :

```
 1   GET /test/_search
 2   {
 3     "query": {
 4       "rank_feature": {
 5         "field": "pagerank",
 6         "saturation": {
 7           "pivot": 8
 8         }
 9       }
10     }
11   }
```

## Relevance Signals

Logarithm – example $\log 4 + pagerank$:

```
GET /test/_search
{
  "query": {
    "rank_feature": {
      "field": "pagerank",
      "log": {
        "scaling_factor": 4
      }
    }
  }
}
```

Sigmoid – example $\frac{pagerank^{0.6}}{pagerank^{0.6}+7^{0.6}}$ :

```
1    GET /test/_search
2    {
3      "query": {
4        "rank_feature": {
5          "field": "pagerank",
6          "sigmoid": {
7            "pivot": 7,
8            "exponent": 0.6
9          }
10        }
11      }
12    }
```

## References

[1] elastic.co. Elasticsearch reference [7.5]: Index moduels:
    Similarity.
    https://www.elastic.co/guide/en/elasticsearch/
    reference/current/index-modules-similarity.html, .
    Accessed: 2020-1-16.

[2] elastic.co. Elasticsearch reference [7.5]: How to: Recipes:
    Incorporating static relevance signals into the score.
    https://www.elastic.co/guide/en/elasticsearch/
    reference/current/static-scoring-signals.html, .
    Accessed: 2020-1-16.

[3] elastic.co. Elasticsearch reference [7.5]: Query dsl: Script score query. https://www.elastic.co/guide/en/elasticsearch/ reference/current/query-dsl-script-score-query.html, . Accessed: 2020-1-16.

[4] elastic.co. Elasticsearch reference [7.5]: Mapping: Similarity. https://www.elastic.co/guide/en/elasticsearch/ reference/current/similarity.html, . Accessed: 2020-1-16.