

Elastic Search: Analyzer[1]

How Text are Tokenized, Selected, and Normalized

Weimao Ke

Drexel University

Table of contents

1. Text Analyzer
2. Index Time Analyzer
3. Search Time Analyzer

Text Analyzer

What is an Analyzer

An Analyzer:

- Converts text into tokens or terms
- May apply filtering, e.g. stopwords removal
- May perform normalization, e.g. case folding
- Built-in or custom

What is an Analyzer

Two types of text analyzers:

- Index time analyzer: converts document text to terms in the inverted index
- Search time analyzer: converts query text to terms for matching

Note:

- Query terms and document terms should be in the same form
- Ideally, use the same index time analyzer and search time analyzer

Index Time Analyzer

Default Analyzer

If the analyzer is not specified before indexing a document:

```
1 PUT /my_index/_doc/1
2 {
3   "message": "The QUICK brown foxes
4               jumped over the lazy dog!"
5 }
```

The default *english* analyzer will be used, which:

- Lowercases each token, e.g. *QUICK* → *quick*
- Removes frequent stopwords, e.g. *the*
- Reduces each term to its word stem, e.g. *foxes* → *fox*

And the following will be used in the inverted index:

```
1 [ quick, brown, fox, jump, over, lazi, dog ]
```

Analyzer Specification

You can specify the index time analyzer with mapping:

```
1  PUT /my_index
2  {
3    "mappings": {
4      "properties": {
5        "message": {
6          "type": "text",
7          "analyzer": "standard"
8        }
9      }
10    }
11  }
```


Analyzer Specification

This line:

```
7  "analyzer": "standard"
```

specifies the *standard* analyzer, which:

- Provides a grammar based tokenization
- Is based on Unicode Text Segmentation
- Works well for many languages

Analyzer Specification

Example:

```
1 POST _analyze
2 {
3   "analyzer": "standard",
4   "text": "The 2 QUICK Brown-Foxes jumped
5           over the lazy dog's bone."
6 }
```

The result of the *standard* analyzer:

```
1 [ the, 2, quick, brown, foxes, jumped,
2   over, the, lazy, dog's, bone ]
```

Analyzer Specification

The *standard* analyzer consists of:

- A standard tokenizer (Unicode Text Segmentation)
- Token filters
 - Lower Case Token Filter
 - Stop Token Filter (disabled by default)

It accepts parameters:

- *max_token_length*: a token must be within the limit or be split.
- *stopwords*: pre-define stopwords (e.g. *_english_*) or a custom array/list.
- *stopwords_path*: the path to a file containing stop words.

Analyzer Specification

Example:

```
1  PUT my_index
2  {
3    "settings": {
4      "analysis": {
5        "analyzer": {
6          "my_english_analyzer": {
7            "type": "standard",
8            "max_token_length": 5,
9            "stopwords": "_english_"
10         }
11       }
12     }
13   }
14 }
```

Analyzer Specification

Now if you use the custom *my_english_analyzer*:

```
1 POST my_index/_analyze
2 {
3   "analyzer": "my_english_analyzer",
4   "text": "The 2 QUICK Brown-Foxes jumped over the
           lazy dog's bone."
5 }
```

The following terms will be produced:

```
1 [ 2, quick, brown, foxes, jumpe, d,
2   over, lazy, dog's, bone ]
```

Note that:

- The stopword "the" has been removed.
- The term "jumped" is split into "jumpe" and "d".

Search Time Analyzer

Search Time Analyzer

With a full-text *match* query, for example:

```
1  GET /my_index/_search
2  {
3      "query": {
4          "match" : {
5              "message" : {
6                  "query" : "a quick fox"
7              }
8          }
9      }
10 }
```

The same analyzer (as the index time analyzer) will be used.

Search Time Analyzer

It is possible, though, to use a different search time analyzer.

To specify a search time analyzer within a query:

```
1  GET /my_index/_search
2  {
3      "query": {
4          "match" : {
5              "message" : {
6                  "query" : "a quick fox",
7                  "analyzer": "standard",
8              }
9          }
10     }
11 }
```


Search Time Analyzer

You may also use a specific search time analyzer with mapping:

```
1  PUT /my_index
2  {
3    "mappings": {
4      "properties": {
5        "message": {
6          "type": "text",
7          "analyzer": "standard",
8          "search_analyzer": "english"
9        }
10     }
11  }
12 }
```

References

- [1] elastic.co. Elasticsearch reference [7.5]: Text analysis.
<https://www.elastic.co/guide/en/elasticsearch/reference/7.5/analysis.html>. Accessed: 2020-1-16.