

1. 分词器概述

分词是构建倒排索引的重要一环。分词根据语言环境的不同可以分为英文分词、中文分词等;根据分词实现的不同又分为标准分词器、空格分词器、停用词分词器等。在传统的分词器不能解决特定业务场景的问题时,往往需要自定义分词器。

1.1 认识分词

对于分词操作来说,英语单词相对而言是更容易辨认和区分的,因为单词之间都会以空格或者标点隔开,举例如下:

```
1 you cannot use from and size to page through more than 10,000 hits
2 you / cannot / use / from / and / size / to / page / through / more / than / 10,000 / hits
```

而中文在单词、句子甚至段落之间没有空格。有些词可以用几个字来表达,但是同样的字在另外的句子中可以拆解成不同的组合。例如

```
1
2 杭州市长春药店
3 杭州 / 市长 / 春药 / 店 (错误)
4 杭州市 / 长春 / 药店 (正确)
```

1.2 为什么需要分词

中文分词是自然语言处理的基础。搜索引擎之所以需要进行中文分词,主要有如下3个维度的原因:

1. **语义维度**: 单字很多时候表达不了语义,而词往往能表达。分词相当于预处理,能使后面和语义有关的分析更准确:
2. **存储维度**: 如果所有文章按照单字来索引,那么所需要的存储空间和搜索计算时间就要多得多:
3. **时间维度**: 通过倒排索引,我们能以 $O(1)$ 的时间复杂度,通过词组找到对应的文章。

以“深入浅出Elasticsearch”这一字符串的检索为例。“深”“入”“浅”“出”这些字在全体内容中可能会无数次出现，如果以这些单独的字为索引，那么就需要添加无数条记录。而以“深入”为索引，所需记录就少了一些；以“深入浅出”为索引，则少得更多；最后以“深入浅出Elasticsearch”为索引，可能就剩余寥寥几条数据。但只有剩余的这些全字符匹配的文档才是我们期望召回的结果。

注意：设计索引的Mapping阶段，要根据业务用途确定是否需要分词。如果不需要分词，则建议设置keyword类型；如果需要分词，则建议设置为text类型并指定分词器。

1.3 分词发生的阶段

写入数据阶段

分词发生在数据写入阶段，也就是数据索引化阶段，其分词逻辑取决于映射参数analyzer。例如，当使用ik_smart分词器对“昨天，小明和他的朋友们去了市中心的图书馆”进行分词后，会将这句话分成不同的词汇或词组。

```
1 POST _analyze
2 {
3   "analyzer": "ik_max_word",
4   "text": "昨天, 小明和他的朋友们去了市中心的图书馆"
5 }
6
7 返回结果:
8 {
9   "tokens": [
10    {
11      "token": "昨天",
12      "start_offset": 0,
13      "end_offset": 2,
14      "type": "CN_WORD",
15      "position": 0
16    },
17    {
18      "token": "小明",
19      "start_offset": 3,
20      "end_offset": 5,
21      "type": "CN_WORD",
22      "position": 1
23    },
24    {
25      "token": "和他",
26      "start_offset": 5,
27      "end_offset": 7,
28      "type": "CN_WORD",
29      "position": 2
30    },
31    {
32      "token": "的",
33      "start_offset": 7,
34      "end_offset": 8,
35      "type": "CN_CHAR",
36      "position": 3
37    },
38    {
39      "token": "朋友们",
```

```
40     "start_offset": 8,
41     "end_offset": 11,
42     "type": "CN_WORD",
43     "position": 4
44 },
45 {
46     "token": "去了",
47     "start_offset": 11,
48     "end_offset": 13,
49     "type": "CN_WORD",
50     "position": 5
51 },
52 {
53     "token": "市中心",
54     "start_offset": 13,
55     "end_offset": 16,
56     "type": "CN_WORD",
57     "position": 6
58 },
59 {
60     "token": "的",
61     "start_offset": 16,
62     "end_offset": 17,
63     "type": "CN_CHAR",
64     "position": 7
65 },
66 {
67     "token": "图书馆",
68     "start_offset": 17,
69     "end_offset": 20,
70     "type": "CN_WORD",
71     "position": 8
72 }
73 ]
74 }
75
```

执行检索阶段

搜索发生时期，其分词仅对搜索词产生作用。在执行“图书馆”检索时，Elasticsearch会根据倒排索引查找所有包含“图书馆”的文档。

2. 分词器的组成

文档被写入并转换为倒排索引之前，Elasticsearch对文档的操作称为分析。而分析是基于Elasticsearch内置分词器(analyzer)或者自定义分词器实现的。分词器由如下三部分组成，如下图所示：

字符过滤器Character Filter

字符过滤器(character filter)将原始文本作为字符流接收，并通过添加、删除或更改字符来转换字符流。

作用：分词之前的预处理，过滤无用字符。

字符过滤器分类如下：

1) HTML Strip Character Filter：用于删除HTML元素，如删除 < b > 标签；解码HTML实体，如将 &转义为&。

```

1 PUT test_html_strip_filter
2 {
3   "settings": {
4     "analysis": {
5       "char_filter": {
6         "my_char_filter": {
7           "type": "html_strip", // html_strip 代表使用 HTML 标签过滤器
8           "escaped_tags": [    // 当前仅保留 a 标签, escaped_tags: 需要保留的 html 标签
9
10              "a"
11            ]
12          }
13        }
14      }
15    }
16 GET test_html_strip_filter/_analyze
17 {
18   "tokenizer": "standard",
19   "char_filter": ["my_char_filter"],
20   "text": ["<p>I&apos;m so <a>happy</a>!</p>"]
21 }

```

2) Mapping Character Filter: 用于替换指定的字符。

```
1 PUT test_html_strip_filter
2 {
3   "settings": {
4     "analysis": {
5       "char_filter": {
6         "my_char_filter": {
7           "type": "mapping", // mapping 代表使用字符映射过滤器
8           "mappings": [ // 数组中规定的字符会被等价替换为 => 指定的字符
9             "滚 => *",
10            "垃 => *",
11            "圾 => *"
12          ]
13        }
14      }
15    }
16  }
17 }
18 GET test_html_strip_filter/_analyze
19 {
20   //"tokenizer": "standard",
21   "char_filter": ["my_char_filter"],
22   "text": "你就是个垃圾！滚"
23 }
```

3) Pattern Replace Character Filter: 可以基于正则表达式替换指定的字符。

```
1 PUT text_pattern_replace_filter
2 {
3   "settings": {
4     "analysis": {
5       "char_filter": {
6         "my_char_filter": {
7           "type": "pattern_replace", // pattern_replace 代表使用正则替换过滤器
8
9           "pattern": "\"\"(\d{3})\d{4}(\d{4})\"\"", // 正则表达式
10          "replacement": "$1****$2"
11        }
12      }
13    }
14  }
15 GET text_pattern_replace_filter/_analyze
16 {
17   "char_filter": ["my_char_filter"],
18   "text": "您的手机号是18868686688"
19 }
```

切词器Tokenizer

若进行了字符过滤，则系统将接收过滤后的字符流；若未进行过滤，则系统接收原始字符流。在接收字符流后，系统将其进行分词，并记录分词后的顺序或位置(position)、起始值(start_offset)以及偏移量(end_offset-start_offset)。而tokenizer负责初步进行文本分词。官方内置了很多种切词器，默认的切词器为 standard。

词项过滤器Token Filter

词项过滤器用来处理切词完成之后的词项，例如把大小写转换，删除停用词或同义词处理等。官方同样预置了很多词项过滤器，基本可以满足日常开发的需要。当然也是支持第三方也自行开发的。


```
1 GET _analyze
2 {
3   "tokenizer" : "standard",
4   "filter" : ["uppercase"],
5   "text" : ["www.elastic.org.cn","www elastic org cn"]
6 }
```

停用词

在切词完成之后，会被干掉词项，即停用词。停用词可以自定义

英文停用词 (**english**) : a, an, and, are, as, at, be, but, by, for, if, in, into, is, it, no, not, of, on, or, such, that, the, their, then, there, these, they, this, to, was, will, with。

中日韩停用词 (**cjk**) : a, and, are, as, at, be, but, by, for, if, in, into, is, it, no, not, of, on, or, s, such, t, that, the, their, then, there, these, they, this, to, was, will, with, www。

```

1 GET _analyze
2 {
3   "tokenizer": "standard",
4   "filter": ["stop"],
5   "text": ["What are you doing"]
6 }
7
8 ### 自定义 filter
9 DELETE test_token_filter_stop
10 PUT test_token_filter_stop
11 {
12   "settings": {
13     "analysis": {
14       "filter": {
15         "my_filter": {
16           "type": "stop",
17           "stopwords": [
18             "www"
19           ],
20           "ignore_case": true
21         }
22       }
23     }
24   }
25 }
26 GET test_token_filter_stop/_analyze
27 {
28   "tokenizer": "standard",
29   "filter": ["my_filter"],
30   "text": ["What www WWW are you doing"]
31 }

```

同义词

同义词定义规则

- a, b, c => d: 这种方式, a、b、c 会被 d 代替。
- a, b, c, d: 这种方式下, a、b、c、d 是等价的。

```
1 PUT test_token_filter_synonym
2 {
3   "settings": {
4     "analysis": {
5       "filter": {
6         "my_synonym": {
7           "type": "synonym",
8           "synonyms": [ "good, nice => excellent" ] //good, nice, excellent
9         }
10      }
11    }
12  }
13 }
14 GET test_token_filter_synonym/_analyze
15 {
16   "tokenizer": "standard",
17   "filter": [ "my_synonym" ],
18   "text": [ "good" ]
19 }
```

实践练习：自定义分词器实现对书籍作者的精确匹配

业务需求是这样的：有一个作者字段，比如Li, LeiLei; Han, MeiMei以及LeiLei Li;现在要对其进行精确匹配。

```
1 POST /booksdemo/_bulk
2 {"index":{"_id":1}}
3 {"name":"Li,LeiLei;Han,MeiMei"}
4 {"index":{"_id":2}}
5 {"name": "LeiLei,Li;MeiMei,Han"}
6
7 # 查不出数据
8 POST /booksdemo/_search
9 {
10   "query": {
11     "match": {
12       "name": "lileilei"
13     }
14   }
15 }
16
```

自定义分词器

```
1 DELETE /booksdemo
2 PUT /booksdemo
3 {
4   "settings": {
5     "analysis": {
6       "char_filter": {
7         "my_char_filter": {
8           "type": "mapping",
9           "mappings": [          //将“,”过滤掉
10             ", => "
11           ]
12         }
13       },
14       "tokenizer": {
15         "my_tokenizer": {
16           "type": "pattern",
17           "pattern": "\"\"\\;\""      //将“;”作为自定义分词分隔符
18         }
19       },
20       "filter": {
21         "my_synonym_filter": {
22           "type": "synonym",
23           "expand": true,
24           "synonyms": [            //添加同义词词组
25             "leileili => lileilei",
26             "meimeihan => hanmeimei"
27           ]
28         }
29       },
30       "analyzer": {
31         "my_analyzer": {
32           "tokenizer": "my_tokenizer",
33           "char_filter": [
34             "my_char_filter"
35           ],
36           "filter": [
37             "lowercase",
38             "my_synonym_filter"
39           ]
40         }
41       }
42     }
43   }
44 }
```

```
40     }
41   }
42 }
43 },
44 "mappings": {
45   "properties": {
46     "name": {
47       "type": "text",
48       "analyzer": "my_analyzer"
49     }
50   }
51 }
52 }
53
54
```

测试自定义分词器效果

```
1 #借助analyzer API验证分词结果是否正确
2 POST booksdemo/_analyze
3 {
4   "analyzer": "my_analyzer",
5   "text": "Li,LeiLei;Han,MeiMei"
6 }
7
8 POST booksdemo/_analyze
9 {
10  "analyzer": "my_analyzer",
11  "text": "LeiLei,Li;MeiMei,Han"
12 }
13
14 POST /booksdemo/_bulk
15 {"index":{"_id":1}}
16 {"name":"Li,LeiLei;Han,MeiMei"}
17 {"index":{"_id":2}}
18 {"name": "LeiLei,Li;MeiMei,Han"}
19
20 POST /booksdemo/_search
21 {
22   "query": {
23     "match": {
24       "name": "lileilei"
25     }
26   }
27 }
```

3. Ngram自定义分词实战

需求背景

当对keyword类型的字段进行高亮查询时，若值为123asd456，查询sd4，则高亮结果是 123asd456。那么，有没有办法只对sd4高亮呢？

用一句话来概括问题：**明明只想查询ID的一部分，但高亮结果是整个ID串，此时应该怎么办？**

解决方案分析


```
1  ###定义索引
2  PUT my_index_phone
3  {
4      "mappings": {
5          "properties": {
6              "phoneNum": {
7                  "type": "keyword"
8              }
9          }
10     }
11 }
12
13 #####批量写入数据
14 POST my_index_phone/_bulk
15 {"index":{"_id":1}}
16 {"phoneNum":"13611112222"}
17 {"index":{"_id":2}}
18 {"phoneNum":"13944248474"}
19
20
21
22
23 ###执行模糊检索和高亮显示
24 POST my_index_phone/_search
25 {
26     "highlight": {
27         "fields": {
28             "phoneNum": {}
29         }
30     },
31     "query": {
32         "bool": {
33             "should": [
34                 {
35                     "wildcard": {
36                         "phoneNum": "*1111*"
37                     }
38                 }
39             ]
40         }
41     }
42 }
```

```
40     }  
41   }  
42 }
```

高亮检索结果如下：

也就是说，整个字符串都呈现为高亮状态了，没有达到预期。检索过程中选择使用wildcard是为了解决子串匹配的问题，wildcard的实现逻辑类似于MySQL的like模糊匹配。传统的text标准分词器，包括中文分词器ik、英文分词器english、standard等都不能解决上述子串匹配问题。

而实际业务需求是这样的：**一方面要求输入子串能召回全串；另一方面要求检索的子串实现高亮。**对此，只能更换一种分词来实现，即Ngram。

Ngram分词实战

Ngram分词定义

Ngram是一种基于统计语言模型的算法。Ngram基本思想是将文本里面的内容按照字节大小进行滑动窗口操作，形成长度是N的字节片段序列。此时每一个字节片段称为gram。对所有gram的出现频度进行统计，并且按照事先设定好的阈值进行过滤，形成关键gram列表，也就是这个文本的向量特征空间。列表中的每一种gram就是一个特征向量维度。该模型基于这样一种假设，第N个词的出现只与前面N-1个词相关，而与其他任何词都不相关，整句的概率就是各个词出现概率的乘积。这些概率可以通过直接从语料中统计N个词同时出现的次数得到。常用的是二元的Bi-Gram（二元语法）和三元的Tri-Gram（三元语法）。

Ngram分词示例

以“你今天吃饭了吗”这一中文句子为例，它的Bi-Gram分词结果如下：

Ngram分词应用场景

- 场景1：文本压缩、检查拼写错误、加速字符串查找、文献语种识别。
- 场景2：自然语言处理自动化领域得到新的应用。如自动分类、自动索引、超链的自动生成、文献检索、无分隔符语言文本的切分等。
- 场景3：自然语言的自动分类功能。针对Elasticsearch检索，Ngram针对无分隔符语言文本的分词（比如手机号检索），可提高检索效率（相较于wildcard检索和正则匹配检索来说）。

Ngram分词实战

```
1 DELETE my_index_phone
2 ###定义索引
3 PUT my_index_phone
4 {
5     "settings":{
6         "number_of_shards":1,
7         "number_of_replicas":0,
8         "index.max_ngram_diff" : 10,
9         "analysis":{
10             "analyzer":{
11                 "phoneNo_analyzer":{
12                     "tokenizer": "phoneNo_analyzer"
13                 }
14             },
15             "tokenizer":{
16                 "phoneNo_analyzer":{
17                     "type": "ngram",
18                     "min_gram": 4,
19                     "max_gram": 11,
20                     "token_chars": [
21                         "letter","digit"
22                     ]
23                 }
24             }
25         }
26     },
27     "mappings":{
28         "dynamic":"strict",
29         "properties":{
30             "phoneNo":{
31                 "type":"text",
32                 "analyzer": "phoneNo_analyzer"
33             }
34         }
35     }
36 }
37
38
39 #####批量写入数据
```

```
40 POST my_index_phone/_bulk
41 {"index":{"_id":1}}
42 {"phoneNo":"13611112222"}
43 {"index":{"_id":2}}
44 {"phoneNo":"13944248474"}
45
```

如上示例共有3个核心参数。

min_gram：最小字符长度（切分），默认为1。

max_gram：最大字符长度（切分），默认为2。

token_chars：表示生成的分词结果中包含的字符类型，默认是全部类型，而在如上的示例中代表保留数字、字母。若只指定letter分词器，则数字就会被过滤掉，分词结果只剩下串中的字符。

借助analyzer API查看分词结果：

```
1
2 POST my_index_phone/_analyze
3 {
4   "analyzer": "phoneNo_analyzer",
5   "text": "13611112222"
6 }
```

检索及高亮的执行语句如下：

```
1 POST my_index_phone/_search
2 {
3   "highlight": {
4     "fields": {
5       "phoneNo": {}
6     }
7   },
8   "query": {
9     "bool": {
10      "should": [
11        {
12          "match_phrase": {
13            "phoneNo": "1111"
14          }
15        }
16      ]
17    }
18  }
19 }
```

返回结果的片段如下:

可以看出, 此时代码已经能满足检索和高亮的双重需求, 也就是说自定义分词完美地解决了提出的问题。