

## 1. ES分词器详解

### 1.1 基本概念

分词器官方称之为文本分析器, 顾名思义, 是对文本进行分析处理的一种手段, 基本处理逻辑为按照预先制定的分词规则, 把原始文档分割成若干更小粒度的词项, 粒度大小取决于分词器规则。

```
#ES的默认分词设置是standard, 会单字拆分
POST _analyze
{
  "analyzer": "standard",
  "text": "中华人民共和国"
}

#ik_smart: 会做最粗粒度的拆
POST _analyze
{
  "analyzer": "ik_smart",
  "text": "中华人民共和国"
}

#ik_max_word: 会将文本做最细粒度的拆分
POST _analyze
{
  "analyzer": "ik_max_word",
  "text": "中华人民共和国"
}
```

```
1 {
2   "tokens" : [
3     {
4       "token" : "中华人民共和国",
5       "start_offset" : 0,
6       "end_offset" : 7,
7       "type" : "CN_WORD",
8       "position" : 0
9     },
10    {
11      "token" : "中华人民",
12      "start_offset" : 0,
13      "end_offset" : 4,
14      "type" : "CN_WORD",
15      "position" : 1
16    },
17    {
18      "token" : "中华",
19      "start_offset" : 0,
20      "end_offset" : 2,
21      "type" : "CN_WORD",
22      "position" : 2
23    },
24  ]
}
```

### 1.2 分词发生时期

分词器的处理过程发生在 Index Time 和 Search Time 两个时期。

- Index Time: 文档写入并创建倒排索引时期, 其分词逻辑取决于映射参数analyzer。
- Search Time: 搜索发生时期, 其分词仅对搜索词产生作用。

### 1.3 分词器的组成

- 切词器 (Tokenizer) : 用于定义切词 (分词) 逻辑
- 词项过滤器 (Token Filter) : 用于对分词之后的单个词项的处理逻辑
- 字符过滤器 (Character Filter) : 用于处理单个字符

注意:

- 分词器不会对源数据造成任何影响, 分词仅仅是对倒排索引或者搜索词的行为。

### 切词器: Tokenizer

tokenizer 是分词器的核心组成部分之一，其主要作用是分词，或称之为切词。主要用来对原始文本进行细粒度拆分。拆分之后的每一个部分称之为一个 Term，或称之为一个词项。可以把切词器理解为预定义的切词规则。官方内置了很多种切词器，默认的切词器位 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   "filter" : ["lowercase"],
4   "text" : "WWW ELASTIC ORG CN"
5 }
6
7 GET _analyze
8 {
9   "tokenizer" : "standard",
10  "filter" : ["uppercase"],
11  "text" : ["www.elastic.org.cn","www elastic org cn"]
12 }
```

## 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 }

```

## 字符过滤器：Character Filter

分词之前的预处理，过滤无用字符。

```

1  PUT <index_name>
2  {
3    "settings": {
4      "analysis": {
5        "char_filter": {
6          "my_char_filter": {
7            "type": "<char_filter_type>"
8          }
9        }
10     }
11  }
12 }

```

type：使用的字符过滤器类型名称，可配置以下值：

- html\_strip
- mapping
- pattern\_replace

## HTML 标签过滤器：HTML Strip Character Filter

字符过滤器会去除 HTML 标签和转义 HTML 元素，如、&

```

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
9              "a"
10           ]
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 }

```

参数: escaped\_tags: 需要保留的 html 标签

## 字符映射过滤器: Mapping Character Filter

通过定义映射替换为规则, 把特定字符替换为指定字符

```

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 }
```

## 正则替换过滤器：Pattern Replace Character Filter

```
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 }
```

## 1.4 倒排索引的数据结构

当数据写入 ES 时，数据将会通过 分词 被切分为不同的 term，ES 将 term 与其对应的文档列表建立一种映射关系，这种结构就是 倒排索引。如下图所示：

为了进一步提升索引的效率，ES 在 term 的基础上利用 term 的前缀或者后缀构建了 term index, 用于对 term 本身进行索引，ES 实际的索引结构如下图所示：

这样当我们去搜索某个关键词时，ES 首先根据它的前缀或者后缀迅速缩小关键词的在 term dictionary 中的范围，大大减少了磁盘IO的次数。

- 单词词典 (Term Dictionary)：记录所有文档的单词，记录单词到倒排列表的关联关系
  - 常用字典数据结构：<https://www.cnblogs.com/LBSer/p/4119841.html>
- 倒排列表(Posting List)-记录了单词对应的文档结合，由倒排索引项组成
- 倒排索引项(Posting):
  - 文档ID
  - 词频TF-该单词在文档中出现的次数，用于相关性评分
  - 位置(Position)-单词在文档中分词的位置。用于短语搜索 (match phrase query)
  - 偏移(Offset)-记录单词的开始结束位置，实现高亮显示

Elasticsearch 的JSON文档中的每个字段，都有自己的倒排索引。

可以指定对某些字段不做索引：

- 优点：节省存储空间
- 缺点: 字段无法被搜索

## 2. 相关性详解

搜索是用户和搜索引擎的对话，**用户关心的是搜索结果的相关性**

- 是否可以找到所有相关的内容
- 有多少不相关的内容被返回了
- 文档的打分是否合理
- 结合业务需求，平衡结果排名

### 2.1 什么是相关性 (Relevance)

**搜索的相关性算分，描述了一个文档和查询语句匹配的程度。**ES 会对每个匹配查询条件的结果进行算分\_score。**打分的本质是排序，需要把最符合用户需求的文档排在前面。**

如下例子：显而易见，查询JAVA多线程设计模式，文档id为2,3的文档的算分更高

关键词	文档ID
-----	------

JAVA	1,2,3
设计模式	1,2,3,4,5,6
多线程	2,3,7,9

如何衡量相关性：

- Precision(查准率)-尽可能返回较少的无关文档
- Recall(查全率)-尽量返回较多的相关文档
- Ranking -是否能够按照相关度进行排序

## 2.2 相关性算法

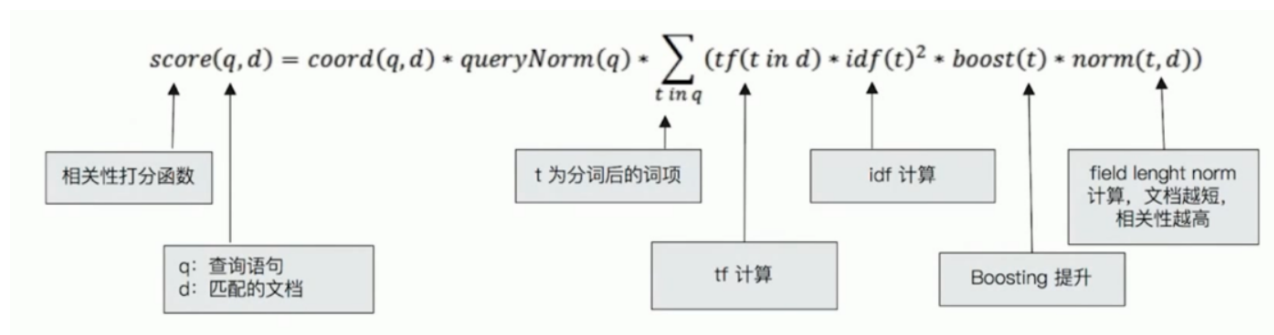
ES 5之前，默认的相关性算分采用TF-IDF，现在采用BM 25。

### TF-IDF

TF-IDF (term frequency-inverse document frequency) 是一种用于信息检索与数据挖掘的常用加权技术。

- TF-IDF被公认为是信息检索领域最重要的发明，除了在信息检索，在文献分类和其他相关领域有着非常广泛的应用。
- IDF的概念，最早是剑桥大学的“斯巴克.琼斯”提出
  - 1972年——“关键词特殊性的统计解释和它在文献检索中的应用”，但是没有从理论上解释IDF应该用  $\log(\text{全部文档数}/\text{检索词出现过的文档总数})$ ，而不是其他函数，也没有做进一步的研究
  - 1970，1980年代萨尔顿和罗宾逊，进行了进一步的证明和研究，并用香农信息论做了证明[http://www.staff.city.ac.uk/~sb317/papers/foundations\\_bm25\\_review.pdf](http://www.staff.city.ac.uk/~sb317/papers/foundations_bm25_review.pdf)
- 现代搜索引擎，对TF-IDF进行了大量细微的优化

Lucene中的TF-IDF评分公式：



- **TF是词频(Term Frequency)**

检索词在文档中出现的频率越高，相关性也越高。



1 词频 (TF) = 某个词在文档中出现的次数 / 文档的总词数

- IDF是逆向文本频率(Inverse Document Frequency)

每个检索词在索引中出现的频率，频率越高，相关性越低。总文档中有些词比如“是”、“的”、“在”在所有文档中出现频率都很高，并不重要，可以减少多个文档中都频繁出现的词的权重。

1 逆向文本频率 (IDF) =  $\log$  (语料库的文档总数 / (包含该词的文档数+1))

- 字段长度归一值 (field-length norm)

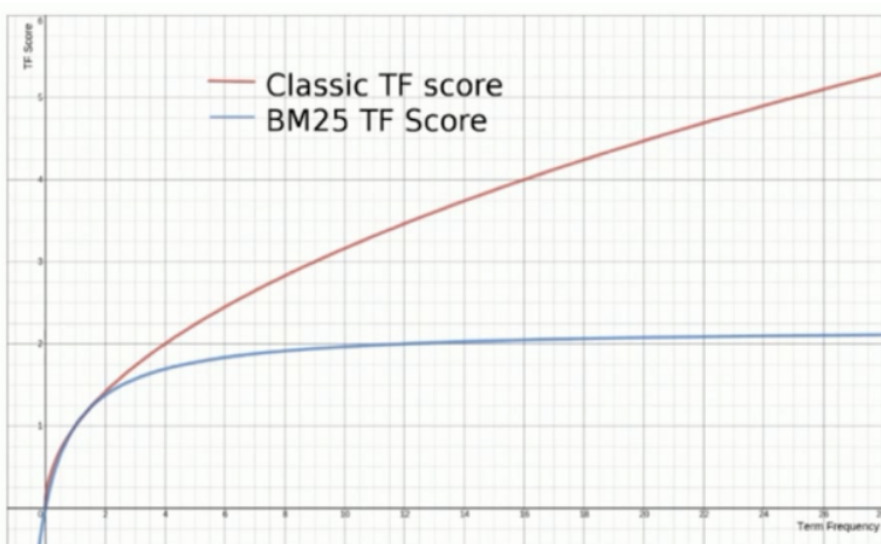
检索词出现在一个内容短的 title 要比同样的词出现在一个内容长的 content 字段权重更大。

以上三个因素——词频 (term frequency)、逆向文本频率 (inverse document frequency) 和字段长度归一值 (field-length norm) ——是在索引时计算并存储的，最后将它们结合在一起计算单个词在特定文档中的权重。

## BM25

BM25 就是对 TF-IDF 算法的改进，对于 TF-IDF 算法，TF(t) 部分的值越大，整个公式返回的值就会越大。BM25 就针对这点进行来优化，随着TF(t) 的逐步加大，该算法的返回值会趋于一个数值。

- 从ES 5开始，默认算法改为BM 25
- 和经典的TF-IDF相比,当TF无限增加时，BM 25算分会趋于一个数值



- BM 25的公式

$$\text{bm25}(d) = \sum_{t \in q, f_{t,d} > 0} \log \left( 1 + \frac{N - df_t + 0.5}{df_t + 0.5} \right) \cdot \frac{f_{t,d}}{f_{t,d} + k \cdot (1 - b + b \frac{1(d)}{\text{avgdl}})}$$

## 2.3 通过Explain API查看TF-IDF

示例:

```
1  PUT /test_score/_bulk
2  {"index":{"_id":1}}
3  {"content":"we use Elasticsearch to power the search"}
4  {"index":{"_id":2}}
5  {"content":"we like elasticsearch"}
6  {"index":{"_id":3}}
7  {"content":"Thre scoring of documents is caculated by the scoring formula"}
8  {"index":{"_id":4}}
9  {"content":"you know,for search"}
10
11 GET /test_score/_search
12 {
13   "explain": true,
14   "query": {
15     "match": {
16       "content": "elasticsearch"
17     }
18   }
19 }
20
21 GET /test_score/_explain/2
22 {
23   "explain": true,
24   "query": {
25     "match": {
26       "content": "elasticsearch"
27     }
28   }
29 }
```

## 2.4 Boosting Query

**Boosting**是控制相关度的一种手段。可以通过指定字段的boost值影响查询结果

## 参数boost的含义:

- 当 $\text{boost} > 1$ 时, 打分的权重相对性提升
- 当 $0 < \text{boost} < 1$ 时, 打分的权重相对性降低
- 当 $\text{boost} < 0$ 时, 贡献负分

应用场景: 希望包含了某项内容的结果不是不出现, 而是排序靠后。

```
1  POST /blogs/_bulk
2  {"index":{"_id":1}}
3  {"title":"Apple iPad","content":"Apple iPad,Apple iPad"}
4  {"index":{"_id":2}}
5  {"title":"Apple iPad,Apple iPad","content":"Apple iPad"}
6
7  GET /blogs/_search
8  {
9    "query": {
10     "bool": {
11       "should": [
12         {
13           "match": {
14             "title": {
15               "query": "apple,ipad",
16               "boost": 1
17             }
18           }
19         },
20         {
21           "match": {
22             "content": {
23               "query": "apple,ipad",
24               "boost": 4
25             }
26           }
27         }
28       ]
29     }
30   }
31 }
```

## 案例：要求苹果公司的产品信息优先展示

```
1 POST /news/_bulk
2 {"index":{"_id":1}}
3 {"content":"Apple Mac"}
4 {"index":{"_id":2}}
5 {"content":"Apple iPad"}
6 {"index":{"_id":3}}
7 {"content":"Apple employee like Apple Pie and Apple Juice"}
8
9
10 GET /news/_search
11 {
12   "query": {
13     "bool": {
14       "must": {
15         "match": {
16           "content": "apple"
17         }
18       }
19     }
20   }
21 }
22
```

## 利用must not排除不是苹果公司产品的文档

```
1 GET /news/_search
2 {
3   "query": {
4     "bool": {
5       "must": {
6         "match": {
7           "content": "apple"
8         }
9       }
10    }
11  }
12 }
```

```
9      },
10      "must_not": {
11        "match": {
12          "content": "pie"
13        }
14      }
15    }
16  }
17 }
```

## 利用negative\_boost降低相关性

对某些返回结果不满意, 但又不想排除掉 ( must\_not), 可以考虑boosting query的 negative\_boost。

- negative\_boost 对 negative部分query生效
- 计算评分时,boosting部分评分不修改, negative部分query乘以negative\_boost值
- negative\_boost取值:0-1.0, 举例:0.3

```
1 GET /news/_search
2 {
3   "query": {
4     "boosting": {
5       "positive": {
6         "match": {
7           "content": "apple"
8         }
9       },
10      "negative": {
11        "match": {
12          "content": "pie"
13        }
14      },
15      "negative_boost": 0.2
16    }
17  }
18 }
```

### 3. 单字符串多字段查询

三种场景：

- **最佳字段(Best Fields)**

当字段之间相互竞争，又相互关联。例如，对于博客的 title 和 body 这样的字段，评分来自最匹配字段

- **多数字段(Most Fields)**

处理英文内容时的一种常见的手段是，在主字段( English Analyzer)，抽取词干，加入同义词，以匹配更多的文档。相同的文本，加入子字段（Standard Analyzer），以提供更加精确的匹配。其他字段作为匹配文档提高相关度的信号，匹配字段越多则越好。

- **混合字段(Cross Field)**

对于某些实体，例如人名，地址，图书信息。需要在多个字段中确定信息，单个字段只能作为整体的一部分。希望在任何这些列出的字段中找到尽可能多的词。

#### 3.1 最佳字段查询Dis Max Query

将任何与任一查询匹配的文档作为结果返回，采用字段上最匹配的评分最终评分返回。

官方文档：<https://www.elastic.co/guide/en/elasticsearch/reference/7.17/query-dsl-dis-max-query.html>

测试

```
1
2 DELETE /blogs
3 PUT /blogs/_doc/1
4 {
5     "title": "Quick brown rabbits",
6     "body": "Brown rabbits are commonly seen."
7 }
8
9 PUT /blogs/_doc/2
10 {
11     "title": "Keeping pets healthy",
12     "body": "My quick brown fox eats rabbits on a regular basis."
13 }
14
```

```

15 POST /blogs/_search
16 {
17     "query": {
18         "bool": {
19             "should": [
20                 { "match": { "title": "Brown fox" } },
21                 { "match": { "body": "Brown fox" } }
22             ]
23         }
24     }
25 }
26

```

思考：查询结果不符合预期，为什么？

```

PUT /blogs/_doc/1
{
  "title": "Quick brown rabbits",
  "body": "Brown rabbits are commonly seen."
}

PUT /blogs/_doc/2
{
  "title": "Keeping pets healthy",
  "body": "My quick brown fox eats rabbits on a regular basis."
}

POST /blogs/_search
{
  "query": {
    "bool": {
      "should": [
        { "match": { "title": "Brown fox" } },
        { "match": { "body": "Brown fox" } }
      ]
    }
  }
}

//
{
  "relation": "eq",
  "max_score": 0.90425634,
  "hits": [
    {
      "_index": "blogs",
      "_type": "_doc",
      "_id": "1",
      "_score": 0.90425634,
      "_source": {
        "title": "Quick brown rabbits",
        "body": "Brown rabbits are commonly seen."
      }
    },
    {
      "_index": "blogs",
      "_type": "_doc",
      "_id": "2",
      "_score": 0.77041256,
      "_source": {
        "title": "Keeping pets healthy",
        "body": "My quick brown fox eats rabbits on a regular basis."
      }
    }
  ]
}

```

## bool should的算法过程：

- 查询should语句中的两个查询
- 加和两个查询的评分
- 乘以匹配语句的总数
- 除以所有语句的总数

上述例子中，title和body属于竞争关系，不应该将分数简单叠加，而是应该找到单个最佳匹配的字段的评分。

## 使用最佳字段查询dis max query

```

1 POST /blogs/_search
2 {
3     "query": {
4         "dis_max": {
5             "queries": [

```

```

6         { "match": { "title": "Brown fox" }},
7         { "match": { "body": "Brown fox" }}
8     ]
9 }
10 }
11 }

```

## 可以通过tie\_breaker参数调整

Tier Breaker是一个介于0-1之间的浮点数。0代表使用最佳匹配;1代表所有语句同等重要。

1. 获得最佳匹配语句的评分\_score。
2. 将其他匹配语句的评分与tie\_breaker相乘
3. 对以上评分求和并规范化

```

1  POST /blogs/_search
2  {
3      "query": {
4          "dis_max": {
5              "queries": [
6                  { "match": { "title": "Quick pets" }},
7                  { "match": { "body": "Quick pets" }}
8              ]
9          }
10     }
11 }
12
13
14 POST /blogs/_search
15 {
16     "query": {
17         "dis_max": {
18             "queries": [
19                 { "match": { "title": "Quick pets" }},
20                 { "match": { "body": "Quick pets" }}
21             ],
22             "tie_breaker": 0.1
23         }
24     }

```



## 3.2 Multi Match Query

### 最佳字段(Best Fields)搜索

**best\_fields**策略获取最佳匹配字段的得分,  $\text{final\_score} = \max(\text{其他匹配字段得分}, \text{最佳匹配字段得分})$

采用 **best\_fields** 查询, 并添加参数  $\text{tie\_breaker}=0.1$ ,  $\text{final\_score} = \text{其他匹配字段得分} * 0.1 + \text{最佳匹配字段得分}$

**Best Fields**是默认类型, 可以不用指定, 等价于**dis\_max**查询方式

```
1 POST /blogs/_search
2 {
3   "query": {
4     "multi_match": {
5       "type": "best_fields",
6       "query": "Brown fox",
7       "fields": ["title", "body"],
8       "tie_breaker": 0.2
9     }
10  }
11 }
```

### 案例

```
1 PUT /employee
2 {
3   "settings" : {
4     "index" : {
5       "analysis.analyzer.default.type": "ik_max_word"
6     }
7   }
8 }
9
10 POST /employee/_bulk
```

```
11 {"index":{"_id":1}}
12 {"empId":"1","name":"员工
001","age":20,"sex":"男","mobile":"19000001111","salary":23343,"deptName":"技术
部","address":"湖北省武汉市洪山区光谷大厦","content":"i like to write best elasticsearch
article"}
13 {"index":{"_id":2}}
14 {"empId":"2","name":"员工
002","age":25,"sex":"男","mobile":"19000002222","salary":15963,"deptName":"销售
部","address":"湖北省武汉市江汉路","content":"i think java is the best programming
language"}
15 {"index":{"_id":3}}
16 {"empId":"3","name":"员工
003","age":30,"sex":"男","mobile":"19000003333","salary":20000,"deptName":"技术
部","address":"湖北省武汉市经济开发区","content":"i am only an elasticsearch beginner"}
17 {"index":{"_id":4}}
18 {"empId":"4","name":"员工
004","age":20,"sex":"女","mobile":"19000004444","salary":15600,"deptName":"销售
部","address":"湖北省武汉市沌口开发区","content":"elasticsearch and hadoop are all very
good solution, i am a beginner"}
19 {"index":{"_id":5}}
20 {"empId":"5","name":"员工
005","age":20,"sex":"男","mobile":"19000005555","salary":19665,"deptName":"测试
部","address":"湖北省武汉市东湖隧道","content":"spark is best big data solution based on
scala, an programming language similar to java"}
21 {"index":{"_id":6}}
22 {"empId":"6","name":"员工
006","age":30,"sex":"女","mobile":"19000006666","salary":30000,"deptName":"技术
部","address":"湖北省武汉市江汉路","content":"i like java developer"}
23 {"index":{"_id":7}}
24 {"empId":"7","name":"员工
007","age":60,"sex":"女","mobile":"19000007777","salary":52130,"deptName":"测试
部","address":"湖北省黄冈市边城区","content":"i like elasticsearch developer"}
25 {"index":{"_id":8}}
26 {"empId":"8","name":"员工
008","age":19,"sex":"女","mobile":"19000008888","salary":60000,"deptName":"技术
部","address":"湖北省武汉市江汉大学","content":"i like spark language"}
27 {"index":{"_id":9}}
28 {"empId":"9","name":"员工
009","age":40,"sex":"男","mobile":"19000009999","salary":23000,"deptName":"销售
部","address":"河南省郑州市郑州大学","content":"i like java developer"}
29 {"index":{"_id":10}}
30 {"empId":"10","name":"张湖
北","age":35,"sex":"男","mobile":"19000001010","salary":18000,"deptName":"测试
部","address":"湖北省武汉市东湖高新","content":"i like java developer, i also like
elasticsearch"}
31 {"index":{"_id":11}}
32 {"empId":"11","name":"王河
南","age":61,"sex":"男","mobile":"19000001011","salary":10000,"deptName":"销售
```

```
部","address":"河南省开封市河南大学","content":"i am not like java"}
33 {"index":{"_id":12}}
34 {"empId":"12","name":"张大
学","age":26,"sex":"女","mobile":"19000001012","salary":11321,"deptName":"测试
部","address":"河南省开封市河南大学","content":"i am java developer, java is good"}
35 {"index":{"_id":13}}
36 {"empId":"13","name":"李江
汉","age":36,"sex":"男","mobile":"19000001013","salary":11215,"deptName":"销售
部","address":"河南省郑州市二七区","content":"i like java and java is very best, i like
it, do you like java"}
37 {"index":{"_id":14}}
38 {"empId":"14","name":"王技
术","age":45,"sex":"女","mobile":"19000001014","salary":16222,"deptName":"测试
部","address":"河南省郑州市金水区","content":"i like c++"}
39 {"index":{"_id":15}}
40 {"empId":"15","name":"张测
试","age":18,"sex":"男","mobile":"19000001015","salary":20000,"deptName":"技术
部","address":"河南省郑州市高新开发区","content":"i think spark is good"}
41
42
43 GET /employee/_search
44 {
45   "query": {
46     "multi_match": {
47       "query": "elasticsearch beginner 湖北省 开封市",
48       "type": "best_fields",
49       "fields": [
50         "content",
51         "address"
52       ]
53     }
54   },
55   "size": 15
56 }
57
58
59 # 查看执行计划
60 GET /employee/_explain/3
61 {
62
63   "query": {
64     "multi_match": {
```

```

65     "query": "elasticsearch beginner 湖北省 开封市",
66     "type": "best_fields",
67     "fields": [
68         "content",
69         "address"
70     ]
71 }
72 }
73 }
74
75 GET /employee/_explain/3
76 {
77
78     "query": {
79         "multi_match": {
80             "query": "elasticsearch beginner 湖北省 开封市",
81             "type": "best_fields",
82             "fields": [
83                 "content",
84                 "address"
85             ],
86             "tie_breaker": 0.1
87         }
88     }
89 }
90

```

### 使用多数字段 (Most Fields) 搜索

**most\_fields**策略获取全部匹配字段的累计得分（综合全部匹配字段的得分），等价于bool should查询方式

```

1 GET /employee/_explain/3
2 {
3
4     "query": {

```

```
5     "multi_match": {
6       "query": "elasticsearch beginner 湖北省 开封市",
7       "type": "most_fields",
8       "fields": [
9         "content",
10        "address"
11      ]
12    }
13  }
14 }
```

## 案例

```
1  DELETE /titles
2  PUT /titles
3  {
4    "mappings": {
5      "properties": {
6        "title": {
7          "type": "text",
8          "analyzer": "english",
9          "fields": {
10            "std": {
11              "type": "text",
12              "analyzer": "standard"
13            }
14          }
15        }
16      }
17    }
18  }
19
20  POST titles/_bulk
21  { "index": { "_id": 1 }}
22  { "title": "My dog barks" }
23  { "index": { "_id": 2 }}
24  { "title": "I see a lot of barking dogs on the road " }
25
```

```
26 # 结果与预期不匹配
27 GET /titles/_search
28 {
29   "query": {
30     "match": {
31       "title": "barking dogs"
32     }
33   }
34 }
```

用广度匹配字段title包括尽可能多的文档——以提升召回率——同时又使用字段title.std 作为信号将相关度更高的文档置于结果顶部。

```
1 GET /titles/_search
2 {
3   "query": {
4     "multi_match": {
5       "query": "barking dogs",
6       "type": "most_fields",
7       "fields": [
8         "title",
9         "title.std"
10      ]
11     }
12   }
13 }
```

每个字段对于最终评分的贡献可以通过自定义值boost 来控制。比如，使title 字段更为重要,这样同时也降低了其他信号字段的作用：

```
1 #增加title的权重
2 GET /titles/_search
3 {
4   "query": {
5     "multi_match": {
6       "query": "barking dogs",
7       "type": "most_fields",
```

```
8     "fields": [  
9         "title^10",  
10        "title.std"  
11    ]  
12 }  
13 }  
14 }
```

### 3.3 跨字段（Cross Field）搜索

搜索内容在多个字段中都显示，类似bool+dis\_max组合

```
1  DELETE /address  
2  PUT /address  
3  {  
4      "settings" : {  
5          "index" : {  
6              "analysis.analyzer.default.type": "ik_max_word"  
7          }  
8      }  
9  }  
10  
11 PUT /address/_bulk  
12 { "index": { "_id": "1" } }  
13 {"province": "湖南", "city": "长沙"}  
14 { "index": { "_id": "2" } }  
15 {"province": "湖南", "city": "常德"}  
16 { "index": { "_id": "3" } }  
17 {"province": "广东", "city": "广州"}  
18 { "index": { "_id": "4" } }  
19 {"province": "湖南", "city": "邵阳"}  
20  
21 #使用most_fields的方式结果不符合预期，不支持operator  
22 GET /address/_search  
23 {  
24     "query": {  
25         "multi_match": {  
26             "query": "湖南常德",
```

```

27     "type": "most_fields",
28     "fields": ["province","city"]
29 }
30 }
31 }
32
33 # 可以使用cross_fields，支持operator
34 #与copy_to相比，其中一个优势就是它可以在搜索时为单个字段提升权重。
35 GET /address/_search
36 {
37     "query": {
38         "multi_match": {
39             "query": "湖南常德",
40             "type": "cross_fields",
41             "operator": "and",
42             "fields": ["province","city"]
43         }
44     }
45 }

```

可以用copy...to 解决，但是需要额外的存储空间

```

1  DELETE /address
2  # copy_to参数允许将多个字段的值复制到组字段中，然后可以将其作为单个字段进行查询
3  PUT /address
4  {
5      "mappings" : {
6          "properties" : {
7              "province" : {
8                  "type" : "keyword",
9                  "copy_to": "full_address"
10             },
11             "city" : {
12                 "type" : "text",
13                 "copy_to": "full_address"
14             }
15         }
16     },

```



```
17     "settings" : {
18         "index" : {
19             "analysis.analyzer.default.type": "ik_max_word"
20         }
21     }
22 }
23
24 PUT /address/_bulk
25 { "index": { "_id": "1" } }
26 {"province": "湖南", "city": "长沙"}
27 { "index": { "_id": "2" } }
28 {"province": "湖南", "city": "常德"}
29 { "index": { "_id": "3" } }
30 {"province": "广东", "city": "广州"}
31 { "index": { "_id": "4" } }
32 {"province": "湖南", "city": "邵阳"}
33
34 GET /address/_search
35 {
36     "query": {
37         "match": {
38             "full_address": {
39                 "query": "湖南常德",
40                 "operator": "and"
41             }
42         }
43     }
44 }
45
```

## 4. Elasticsearch聚合操作

Elasticsearch除搜索以外，提供了针对ES 数据进行统计分析的功能。聚合(aggregations)可以让我们极其方便的实现对数据的统计、分析、运算。例如：

- 什么品牌的手机最受欢迎？
- 这些手机的平均价格、最高价格、最低价格？

- 这些手机每月的销售情况如何？

## 使用场景

聚合查询可以用于各种场景，比如商业智能、数据挖掘、日志分析等等。

- 电商平台的销售分析：统计每个地区的销售额、每个用户的消费总额、每个产品的销售量等，以便更好地了解销售情况和趋势。
- 社交媒体的用户行为分析：统计每个用户的发布次数、转发次数、评论次数等，以便更好地了解用户行为和趋势，同时可以将数据按照地区、时间、话题等维度进行分析。
- 物流企业的运输分析：统计每个区域的运输量、每个车辆的运输次数、每个司机的行驶里程等，以便更好地了解运输情况和优化运输效率。
- 金融企业的交易分析：统计每个客户的交易总额、每个产品的销售量、每个交易员的业绩等，以便更好地了解交易情况和优化业务流程。
- 智能家居的设备监控分析：统计每个设备的使用次数、每个家庭的能源消耗量、每个时间段的设备使用率等，以便更好地了解用户需求和优化设备效能。

## 基本语法

聚合查询的语法结构与其他查询相似，通常包含以下部分：

- 查询条件：指定需要聚合的文档，可以使用标准的 Elasticsearch 查询语法，如 term、match、range 等等。
- 聚合函数：指定要执行的聚合操作，如 sum、avg、min、max、terms、date\_histogram 等等。每个聚合命令都会生成一个聚合结果。
- 聚合嵌套：聚合命令可以嵌套，以便更细粒度地分析数据。

```
1 GET <index_name>/_search
2 {
3   "aggs": {
4     "<aggs_name>": { // 聚合名称需要自己定义
5       "<agg_type>": {
6         "field": "<field_name>"
7       }
8     }
9   }
10 }
```

- aggs\_name：聚合函数的名称
- agg\_type：聚合种类，比如是桶聚合（terms）或者是指标聚合（avg、sum、min、max等）
- field\_name：字段名称或者叫域名。

### 4.1 聚合的分类

- Metric Aggregation: 一些数学运算, 可以对文档字段进行统计分析, 类比Mysql中的 min(), max(), sum() 操作。

```
1 SELECT MIN(price), MAX(price) FROM products
2 #Metric聚合的DSL类比实现:
3 {
4     "aggs":{
5         "avg_price":{
6             "avg":{
7                 "field":"price"
8             }
9         }
10    }
11 }
```

- Bucket Aggregation: 一些满足特定条件的文档的集合放置到一个桶里, 每一个桶关联一个key, 类比Mysql中的group by操作。

```
1 ELECT size COUNT(*) FROM products GROUP BY size
2 #bucket聚合的DSL类比实现:
3 {
4     "aggs": {
5         "by_size": {
6             "terms": {
7                 "field": "size"
8             }
9         }
10    }
```

- Pipeline Aggregation: 对其他的聚合结果进行二次聚合

## 示例数据

```
1 DELETE /employees
2 #创建索引库
3 PUT /employees
```

```
4 {
5   "mappings": {
6     "properties": {
7       "age":{
8         "type": "integer"
9       },
10      "gender":{
11        "type": "keyword"
12      },
13      "job":{
14        "type" : "text",
15        "fields" : {
16          "keyword" : {
17            "type" : "keyword",
18            "ignore_above" : 50
19          }
20        }
21      },
22      "name":{
23        "type": "keyword"
24      },
25      "salary":{
26        "type": "integer"
27      }
28    }
29  }
30 }
31
32 PUT /employees/_bulk
33 { "index" : { "_id" : "1" } }
34 { "name" : "Emma", "age":32, "job":"Product Manager", "gender":"female", "salary":35000 }
35 { "index" : { "_id" : "2" } }
36 { "name" : "Underwood", "age":41, "job":"Dev Manager", "gender":"male", "salary": 50000}
37 { "index" : { "_id" : "3" } }
38 { "name" : "Tran", "age":25, "job":"Web Designer", "gender":"male", "salary":18000 }
39 { "index" : { "_id" : "4" } }
40 { "name" : "Rivera", "age":26, "job":"Web Designer", "gender":"female", "salary": 22000}
41 { "index" : { "_id" : "5" } }
42 { "name" : "Rose", "age":25, "job":"QA", "gender":"female", "salary":18000 }
43 { "index" : { "_id" : "6" } }
```

```

44 { "name" : "Lucy", "age":31, "job":"QA", "gender":"female", "salary": 25000}
45 { "index" : { "_id" : "7" } }
46 { "name" : "Byrd", "age":27, "job":"QA", "gender":"male", "salary":20000 }
47 { "index" : { "_id" : "8" } }
48 { "name" : "Foster", "age":27, "job":"Java Programmer", "gender":"male", "salary": 20000}
49 { "index" : { "_id" : "9" } }
50 { "name" : "Gregory", "age":32, "job":"Java Programmer", "gender":"male", "salary":22000 }
51 { "index" : { "_id" : "10" } }
52 { "name" : "Bryant", "age":20, "job":"Java Programmer", "gender":"male", "salary": 9000}
53 { "index" : { "_id" : "11" } }
54 { "name" : "Jenny", "age":36, "job":"Java Programmer", "gender":"female", "salary":38000 }
55 { "index" : { "_id" : "12" } }
56 { "name" : "Mcdonald", "age":31, "job":"Java Programmer", "gender":"male", "salary": 32000}
57 { "index" : { "_id" : "13" } }
58 { "name" : "Jonthna", "age":30, "job":"Java Programmer", "gender":"female", "salary":30000
    }
59 { "index" : { "_id" : "14" } }
60 { "name" : "Marshall", "age":32, "job":"Javascript Programmer", "gender":"male", "salary":
    25000}
61 { "index" : { "_id" : "15" } }
62 { "name" : "King", "age":33, "job":"Java Programmer", "gender":"male", "salary":28000 }
63 { "index" : { "_id" : "16" } }
64 { "name" : "Mccarthy", "age":21, "job":"Javascript Programmer", "gender":"male", "salary":
    16000}
65 { "index" : { "_id" : "17" } }
66 { "name" : "Goodwin", "age":25, "job":"Javascript Programmer", "gender":"male", "salary":
    16000}
67 { "index" : { "_id" : "18" } }
68 { "name" : "Catherine", "age":29, "job":"Javascript
    Programmer", "gender":"female", "salary": 20000}
69 { "index" : { "_id" : "19" } }
70 { "name" : "Boone", "age":30, "job":"DBA", "gender":"male", "salary": 30000}
71 { "index" : { "_id" : "20" } }
72 { "name" : "Kathy", "age":29, "job":"DBA", "gender":"female", "salary": 20000}

```

## 4.2 Metric Aggregation

- 单值分析：只输出一个分析结果
  - min, max, avg, sum

- Cardinality (类似distinct Count)
- 多值分析:输出多个分析结果
  - stats (统计) , extended stats
  - percentile (百分位) , percentile rank
  - top hits(排在前面的示例)

## 查询员工的最低最高和平均工资

```
1 #多个 Metric 聚合，找到最低最高和平均工资
2 POST /employees/_search
3 {
4   "size": 0,
5   "aggs": {
6     "max_salary": {
7       "max": {
8         "field": "salary"
9       }
10    },
11    "min_salary": {
12      "min": {
13        "field": "salary"
14      }
15    },
16    "avg_salary": {
17      "avg": {
18        "field": "salary"
19      }
20    }
21  }
22 }
```

## 对salary进行统计

```
1 # 一个聚合，输出多值
2 POST /employees/_search
3 {
```

```
4  "size": 0,
5  "aggs": {
6    "stats_salary": {
7      "stats": {
8        "field": "salary"
9      }
10   }
11 }
12 }
```

## cardinate对搜索结果去重

```
1  POST /employees/_search
2  {
3    "size": 0,
4    "aggs": {
5      "cardinate": {
6        "cardinality": {
7          "field": "job.keyword"
8        }
9      }
10   }
11 }
```

## 4.3 Bucket Aggregation

按照一定的规则，将文档分配到不同的桶中，从而达到分类的目的。ES提供的一些常见的 Bucket Aggregation。

- Terms, 需要字段支持filedata
  - keyword 默认支持fielddata
  - text需要在Mapping 中开启fielddata, 会按照分词后的结果进行分桶
- 数字类型

- Range / Data Range
- Histogram (直方图) / Date Histogram
- 支持嵌套: 也就在桶里再做分桶

## 获取job的分类信息

```
1 # 对keyword 进行聚合
2 GET /employees/_search
3 {
4   "size": 0,
5   "aggs": {
6     "jobs": {
7       "terms": {
8         "field": "job.keyword"
9       }
10    }
11  }
```

聚合可配置属性有:

- field: 指定聚合字段
- size: 指定聚合结果数量
- order: 指定聚合结果排序方式

默认情况下, Bucket聚合会统计Bucket内的文档数量, 记为\_count, 并且按照\_count降序排序。我们可以指定order属性, 自定义聚合的排序方式:

```
1 GET /employees/_search
2 {
3   "size": 0,
4   "aggs": {
5     "jobs": {
6       "terms": {
7         "field": "job.keyword",
8         "size": 10,
9         "order": {
10          "_count": "desc"
11        }
12      }
13    }
14  }
```



```
11     }
12   }
13 }
14 }
15 }
```

## 限定聚合范围

```
1  #只对salary在10000元以上的文档聚合
2  GET /employees/_search
3  {
4    "query": {
5      "range": {
6        "salary": {
7          "gte": 10000
8        }
9      }
10 },
11 "size": 0,
12 "aggs": {
13   "jobs": {
14     "terms": {
15       "field": "job.keyword",
16       "size": 10,
17       "order": {
18         "_count": "desc"
19       }
20     }
21   }
22 }
23 }
```

注意：对 Text 字段进行 terms 聚合查询，会失败抛出异常

```
1  POST /employees/_search
```

```
2  {
3    "size": 0,
4    "aggs": {
5      "jobs": {
6        "terms": {
7          "field": "job"
8        }
9      }
10   }
11 }
```

解决办法：对 Text 字段打开 fielddata，支持terms aggregation

```
1  PUT /employees/_mapping
2  {
3    "properties" : {
4      "job":{
5        "type": "text",
6        "fielddata": true
7      }
8    }
9  }
10
11 # 对 Text 字段进行分词，分词后的terms
12 POST /employees/_search
13 {
14   "size": 0,
15   "aggs": {
16     "jobs": {
17       "terms": {
18         "field": "job"
19       }
20     }
21   }
22 }
```

对job.keyword 和 job 进行 terms 聚合，分桶的总数并不一样

```
1 POST /employees/_search
2 {
3   "size": 0,
4   "aggs": {
5     "cardinate": {
6       "cardinality": {
7         "field": "job"
8       }
9     }
10  }
11 }
```

## Range & Histogram聚合

- 按照数字的范围，进行分桶
- 在Range Aggregation中，可以自定义Key

Range 示例：按照工资的 Range 分桶

```
1 Salary Range分桶，可以自己定义 key
2 POST employees/_search
3 {
4   "size": 0,
5   "aggs": {
6     "salary_range": {
7       "range": {
8         "field": "salary",
9         "ranges": [
10          {
11            "to": 10000
12          },
13          {
14            "from": 10000,
```

```

15         "to":20000
16     },
17     {
18         "key": ">20000",
19         "from":20000
20     }
21 ]
22 }
23 }
24 }
25 }

```

## Histogram示例：按照工资的间隔分桶

```

1  #工资0到10万，以 5000一个区间进行分桶
2  POST employees/_search
3  {
4      "size": 0,
5      "aggs": {
6          "salary_histogram": {
7              "histogram": {
8                  "field":"salary",
9                  "interval":5000,
10                 "extended_bounds":{
11                     "min":0,
12                     "max":100000
13                 }
14             }
15         }
16     }
17 }

```

top\_hits应用场景: 当获取分桶后，桶内最匹配的顶部文档列表

```

1 # 指定size, 不同工种中, 年纪最大的3个员工的具体信息
2 POST /employees/_search
3 {
4   "size": 0,
5   "aggs": {
6     "jobs": {
7       "terms": {
8         "field": "job.keyword"
9       },
10      "aggs": {
11        "old_employee": {
12          "top_hits": {
13            "size": 3,
14            "sort": [
15              {
16                "age": {
17                  "order": "desc"
18                }
19              }
20            ]
21          }
22        }
23      }
24    }
25  }
26 }
27

```

## 嵌套聚合示例

```

1 # 嵌套聚合1, 按照工作类型分桶, 并统计工资信息
2 POST employees/_search
3 {
4   "size": 0,
5   "aggs": {

```

```
6     "Job_salary_stats": {
7         "terms": {
8             "field": "job.keyword"
9         },
10        "aggs": {
11            "salary": {
12                "stats": {
13                    "field": "salary"
14                }
15            }
16        }
17    }
18 }
19 }
20
21 # 多次嵌套。根据工作类型分桶，然后按照性别分桶，计算工资的统计信息
22 POST employees/_search
23 {
24     "size": 0,
25     "aggs": {
26         "Job_gender_stats": {
27             "terms": {
28                 "field": "job.keyword"
29             },
30             "aggs": {
31                 "gender_stats": {
32                     "terms": {
33                         "field": "gender"
34                     },
35                     "aggs": {
36                         "salary_stats": {
37                             "stats": {
38                                 "field": "salary"
39                             }
40                         }
41                     }
42                 }
43             }
44         }
45     }
```

## 4.4 Pipeline Aggregation

支持对聚合分析的结果，再次进行聚合分析。

Pipeline 的分析结果会输出到原结果中，根据位置的不同，分为两类：

- Sibling - 结果和现有分析结果同级
  - Max, min, Avg & Sum Bucket
  - Stats, Extended Status Bucket
  - Percentiles Bucket
- Parent - 结果内嵌到现有的聚合分析结果之中
  - Derivative(求导)
  - Cumulative Sum(累计求和)
  - Moving Function(移动平均值)

### min\_bucket示例

在员工数最多的工种里，找出平均工资最低的工种

```

1 # 平均工资最低的工种
2 POST employees/_search
3 {
4   "size": 0,
5   "aggs": {
6     "jobs": {
7       "terms": {
8         "field": "job.keyword",
9         "size": 10
10      },
11     "aggs": {
12       "avg_salary": {
13         "avg": {
14           "field": "salary"
15         }
16       }
17     }
18   }
19 }
```

```

18     },
19     "min_salary_by_job":{
20         "min_bucket": {
21             "buckets_path": "jobs>avg_salary"
22         }
23     }
24 }
25 }

```

- min\_salary\_by\_job结果和jobs的聚合同级
- min\_bucket求之前结果的最小值
- 通过bucket\_path关键字指定路径

## Stats示例

```

1  # 平均工资的统计分析
2  POST employees/_search
3  {
4      "size": 0,
5      "aggs": {
6          "jobs": {
7              "terms": {
8                  "field": "job.keyword",
9                  "size": 10
10             },
11             "aggs": {
12                 "avg_salary": {
13                     "avg": {
14                         "field": "salary"
15                     }
16                 }
17             }
18         },
19         "stats_salary_by_job":{
20             "stats_bucket": {
21                 "buckets_path": "jobs>avg_salary"
22             }
23         }
24     }

```



```
25 }
```

## percentiles示例

```
1 # 平均工资的百分位数
2 POST employees/_search
3 {
4   "size": 0,
5   "aggs": {
6     "jobs": {
7       "terms": {
8         "field": "job.keyword",
9         "size": 10
10      },
11      "aggs": {
12        "avg_salary": {
13          "avg": {
14            "field": "salary"
15          }
16        }
17      }
18    },
19    "percentiles_salary_by_job": {
20      "percentiles_bucket": {
21        "buckets_path": "jobs>avg_salary"
22      }
23    }
24  }
25 }
26
```

## Cumulative\_sum示例

```
1 #Cumulative_sum    累计求和
2 POST employees/_search
```

```

3  {
4    "size": 0,
5    "aggs": {
6      "age": {
7        "histogram": {
8          "field": "age",
9          "min_doc_count": 0,
10         "interval": 1
11       },
12       "aggs": {
13         "avg_salary": {
14           "avg": {
15             "field": "salary"
16           }
17         },
18         "cumulative_salary": {
19           "cumulative_sum": {
20             "buckets_path": "avg_salary"
21           }
22         }
23       }
24     }
25   }
26 }
27

```

## 4.5 聚合的作用范围

ES聚合分析的默认作用范围是query的查询结果集，同时ES还支持以下方式改变聚合的作用范围：

- Filter
- Post Filter
- Global

```

1  #Query
2  POST employees/_search
3  {

```

```
4  "size": 0,
5  "query": {
6    "range": {
7      "age": {
8        "gte": 20
9      }
10   }
11 },
12 "aggs": {
13   "jobs": {
14     "terms": {
15       "field": "job.keyword"
16     }
17   }
18 }
19 }
20 }
```

```
21
22 #Filter
23 POST employees/_search
24 {
25   "size": 0,
26   "aggs": {
27     "older_person": {
28       "filter": {
29         "range": {
30           "age": {
31             "from": 35
32           }
33         }
34       },
35       "aggs": {
36         "jobs": {
37           "terms": {
38             "field": "job.keyword"
39           }
40         }
41       }
42     },
43     "all_jobs": {
44       "terms": {
```

```
44         "field": "job.keyword"
```

```
45
```

```
46     }
```

```
47 }
```

```
48 }
```

```
49 }
```

```
50
```

```
51
```

```
52
```

53 **#Post field**。一条语句，找出所有的**job**类型。还能找到聚合后符合条件的结果

```
54 POST employees/_search
```

```
55 {
```

```
56   "aggs": {
```

```
57     "jobs": {
```

```
58       "terms": {
```

```
59         "field": "job.keyword"
```

```
60       }
```

```
61     }
```

```
62   },
```

```
63   "post_filter": {
```

```
64     "match": {
```

```
65       "job.keyword": "Dev Manager"
```

```
66     }
```

```
67   }
```

```
68 }
```

```
69
```

```
70
```

```
71 #global
```

```
72 POST employees/_search
```

```
73 {
```

```
74   "size": 0,
```

```
75   "query": {
```

```
76     "range": {
```

```
77       "age": {
```

```
78         "gte": 40
```

```
79       }
```

```
80     }
```

```
81   },
```

```
82   "aggs": {
```

```

83     "jobs": {
84         "terms": {
85             "field": "job.keyword"
86
87         }
88     },
89
90     "all": {
91         "global": {},
92         "aggs": {
93             "salary_avg": {
94                 "avg": {
95                     "field": "salary"
96                 }
97             }
98         }
99     }
100 }
101 }
102
103
104

```

## 4.6 排序

指定order，按照count和key进行排序：

- 默认情况，按照count降序排序
- 指定size，就能返回相应的桶

```

1  #排序 order
2  #count and key
3  POST employees/_search
4  {
5      "size": 0,
6      "query": {
7          "range": {
8              "age": {

```

```

9         "gte": 20
10     }
11 }
12 },
13 "aggs": {
14     "jobs": {
15         "terms": {
16             "field": "job.keyword",
17             "order": [
18                 {"_count": "asc"},
19                 {"_key": "desc"}
20             ]
21         }
22     }
23 }
24 }
25 }
26
27
28 #排序 order
29 #count and key
30 POST employees/_search
31 {
32     "size": 0,
33     "aggs": {
34         "jobs": {
35             "terms": {
36                 "field": "job.keyword",
37                 "order": [ {
38                     "avg_salary": "desc"
39                 } ]
40             }
41         },
42     },
43     "aggs": {
44         "avg_salary": {
45             "avg": {
46                 "field": "salary"
47             }

```

```

48     }
49 }
50 }
51 }
52 }
53
54
55 #排序 order
56 #count and key
57 POST employees/_search
58 {
59   "size": 0,
60   "aggs": {
61     "jobs": {
62       "terms": {
63         "field": "job.keyword",
64         "order": [ {
65           "stats_salary.min": "desc"
66         }]
67       }
68     },
69     "stats_salary": {
70       "stats": {
71         "field": "salary"
72       }
73     }
74   }
75 }
76 }
77 }
78 }
79 }

```

## 4.7 ES聚合分析不精准原因分析

ElasticSearch在对海量数据进行聚合分析的时候会损失搜索的精准度来满足实时性的需求。

Terms聚合分析的执行流程：

不精准的原因：数据分散到多个分片，聚合是每个分片的取 Top X，导致结果不精准。ES 可以不每个分片 Top X，而是全量聚合，但势必这会有很大的性能问题。

**思考：如何提高聚合精确度？**

### 方案1：设置主分片为1

注意7.x版本已经默认为1。

适用场景：数据量小的小集群规模业务场景。

### 方案2：调大 shard\_size 值

设置 shard\_size 为比较大的值，官方推荐： $size * 1.5 + 10$ 。shard\_size 值越大，结果越趋近于精准聚合结果值。此外，还可以通过 show\_term\_doc\_count\_error 参数显示最差情况下的错误值，用于辅助确定 shard\_size 大小。

- size: 是聚合结果的返回值，客户期望返回聚合排名前三，size 值就是 3。
- shard\_size: 每个分片上聚合的数据条数。shard\_size 原则上要大于等于 size

适用场景：数据量大、分片数多的集群业务场景。

测试：使用 kibana 的测试数据

```
1 DELETE my_flights
2 PUT my_flights
3 {
4   "settings": {
5     "number_of_shards": 20
6   },
7   "mappings" : {
8     "properties" : {
9       "AvgTicketPrice" : {
10        "type" : "float"
11      },
12      "Cancelled" : {
13        "type" : "boolean"
14      },
15      "Carrier" : {
16        "type" : "keyword"
17      },
18      "Dest" : {
```



```
19     "type" : "keyword"
20 },
21     "DestAirportID" : {
22         "type" : "keyword"
23     },
24     "DestCityName" : {
25         "type" : "keyword"
26     },
27     "DestCountry" : {
28         "type" : "keyword"
29     },
30     "DestLocation" : {
31         "type" : "geo_point"
32     },
33     "DestRegion" : {
34         "type" : "keyword"
35     },
36     "DestWeather" : {
37         "type" : "keyword"
38     },
39     "DistanceKilometers" : {
40         "type" : "float"
41     },
42     "DistanceMiles" : {
43         "type" : "float"
44     },
45     "FlightDelay" : {
46         "type" : "boolean"
47     },
48     "FlightDelayMin" : {
49         "type" : "integer"
50     },
51     "FlightDelayType" : {
52         "type" : "keyword"
53     },
54     "FlightNum" : {
55         "type" : "keyword"
56     },
57     "FlightTimeHour" : {
58         "type" : "keyword"
```

```
59     },
60     "FlightTimeMin" : {
61         "type" : "float"
62     },
63     "Origin" : {
64         "type" : "keyword"
65     },
66     "OriginAirportID" : {
67         "type" : "keyword"
68     },
69     "OriginCityName" : {
70         "type" : "keyword"
71     },
72     "OriginCountry" : {
73         "type" : "keyword"
74     },
75     "OriginLocation" : {
76         "type" : "geo_point"
77     },
78     "OriginRegion" : {
79         "type" : "keyword"
80     },
81     "OriginWeather" : {
82         "type" : "keyword"
83     },
84     "dayOfWeek" : {
85         "type" : "integer"
86     },
87     "timestamp" : {
88         "type" : "date"
89     }
90 }
91 }
92 }
93
94 POST _reindex
95 {
96     "source": {
97         "index": "kibana_sample_data_flights"
```

```

98   },
99   "dest": {
100     "index": "my_flights"
101   }
102 }
103
104 GET my_flights/_count
105 GET kibana_sample_data_flights/_search
106 {
107   "size": 0,
108   "aggs": {
109     "weather": {
110       "terms": {
111         "field": "OriginWeather",
112         "size": 5,
113         "show_term_doc_count_error": true
114       }
115     }
116   }
117 }
118
119 GET my_flights/_search
120 {
121   "size": 0,
122   "aggs": {
123     "weather": {
124       "terms": {
125         "field": "OriginWeather",
126         "size": 5,
127         "shard_size": 10,
128         "show_term_doc_count_error": true
129       }
130     }
131   }
132 }

```

在Terms Aggregation的返回中有两个特殊的数值：

- `doc_count_error_upper_bound` : 被遗漏的term 分桶，包含的文档，有可能的最大值
- `sum_other_doc_count`: 除了返回结果 bucket的terms以外，其他 terms 的文档总数（总数-返回的总数）

### 方案3：将size设置为全量值，来解决精度问题

将size设置为2的32次方减去1也就是分片支持的最大值，来解决精度问题。

原因：1.x版本，size等于0代表全部，高版本取消0值，所以设置了最大值（大于业务的全量值）。全量带来的弊端就是：如果分片数据量极大，这样做会耗费巨大的CPU资源来排序，而且可能会阻塞网络。

适用场景：对聚合精准度要求极高的业务场景，由于性能问题，不推荐使用。

### 方案4：使用Clickhouse/ Spark 进行精准聚合

适用场景：数据量非常大、聚合精度要求高、响应速度快的业务场景。

## 4.8 Elasticsearch 聚合性能优化

### 启用 eager global ordinals 提升高基数聚合性能

适用场景：高基数聚合。高基数聚合场景中的高基数含义：一个字段包含很大比例的唯一值。

global ordinals 中文翻译成全局序号，是一种数据结构，应用场景如下：

- 基于 keyword, ip 等字段的分桶聚合，包含：terms聚合、composite 聚合等。
- 基于text 字段的分桶聚合（前提条件是：fielddata 开启）。
- 基于父子文档 Join 类型的 has\_child 查询和 父聚合。

global ordinals 使用一个数值代表字段中的字符串值，然后为每一个数值分配一个 bucket（分桶）。

global ordinals 的本质是：启用 eager\_global\_ordinals 时，会在刷新（refresh）分片时构建全局序号。这将构建全局序号的成本从搜索阶段转移到了数据索引化（写入）阶段。

创建索引的同时开启：eager\_global\_ordinals。

```
1 PUT /my-index
2 {
3   "mappings": {
4     "properties": {
5       "tags": {
6         "type": "keyword",
7         "eager_global_ordinals": true
8       }
9     }
10  }
```

```
10 }
```

注意：开启 `eager_global_ordinals` 会影响写入性能，因为每次刷新时都会创建新的全局序号。为了最大程度地减少由于频繁刷新建立全局序号而导致的额外开销，请调大刷新间隔 `refresh_interval`。动态调整刷新频率的方法如下：

```
1 PUT my-index/_settings
2 {
3   "index": {
4     "refresh_interval": "30s"
5   }
}
```

该招数的本质是：以空间换时间。

## 插入数据时对索引进行预排序

- **Index sorting**（索引排序）可用于在插入时对索引进行预排序，而不是在查询时再对索引进行排序，这将提高范围查询（range query）和排序操作的性能。
- 在 Elasticsearch 中创建新索引时，可以配置如何对每个分片内的段进行排序。
- 这是 Elasticsearch 6.X 之后版本才有的特性。

```
1
2 PUT /my_index
3 {
4   "settings": {
5     "index": {
6       "sort.field": "create_time",
7       "sort.order": "desc"
8     }
9   },
10  "mappings": {
11    "properties": {
12      "create_time": {
13        "type": "date"
14      }
15    }
16  }
```

```
16     }
17 }
```

**注意：预排序将增加 Elasticsearch 写入的成本。**在某些用户特定场景下，开启索引预排序会导致大约 40%-50% 的写性能下降。也就是说，如果用户场景更关注写性能的业务，开启索引预排序不是一个很好的选择。

## 使用节点查询缓存

**节点查询缓存（Node query cache）**可用于有效缓存过滤器（filter）操作的结果。如果多次执行同一 filter 操作，这将很有效，但是即便更改过滤器中的某一个值，也将意味着需要计算新的过滤器结果。

例如，由于 “now” 值一直在变化，因此无法缓存在过滤器上下文中使用 “now” 的查询。

那怎么使用缓存呢？通过在 now 字段上应用 datemath 格式将其四舍五入到最接近的分钟/小时等，可以使此类请求更具可缓存性，以便可以对筛选结果进行缓存。

```
1  PUT /my_index/_doc/1
2  {
3    "create_time":"2022-05-11T16:30:55.328Z"
4  }
5
6  #下面的示例无法使用缓存
7  GET /my_index/_search
8  {
9    "query":{
10      "constant_score": {
11        "filter": {
12          "range": {
13            "create_time": {
14              "gte": "now-1h",
15              "lte": "now"
16            }
17          }
18        }
19      }
20    }
21  }
22
```

23 # 下面的示例就可以使用节点查询缓存。

24 GET /my\_index/\_search

```
25 {
26   "query":{
27     "constant_score": {
28       "filter": {
29         "range": {
30           "create_time": {
31             "gte": "now-1h/m",
32             "lte": "now/m"
33           }
34         }
35       }
36     }
37   }
38 }
```

上述示例中的“now-1h/m”就是 datemath 的格式。

如果当前时间 now 是：16:31:29，那么range query 将匹配 my\_date 介于：15:31:00 和 15:31:59 之间的时间数据。同理，聚合的前半部分 query 中如果有基于时间查询，或者后半部分 aggs 部分中有基于时间聚合的，建议都使用 datemath 方式做缓存处理以优化性能。

## 使用分片请求缓存

聚合语句中，设置：size：0，就会使用分片请求缓存缓存结果。size = 0 的含义是：只返回聚合结果，不返回查询结果。

```
1 GET /es_db/_search
2 {
3   "size": 0,
4   "aggs": {
5     "remark_agg": {
6       "terms": {
7         "field": "remark.keyword"
8       }
9     }
10  }
11 }
```

## 拆分聚合，使聚合并行化

Elasticsearch 查询条件中同时有多个条件聚合，默认情况下聚合不是并行运行的。当为每个聚合提供自己的查询并执行 `msearch` 时，性能会有显著提升。因此，在 CPU 资源不是瓶颈的前提下，如果想缩短响应时间，可以将多个聚合拆分为多个查询，借助：[msearch 实现并行聚合](#)。

```
1 #常规的多条件聚合实现
2 GET /employees/_search
3 {
4   "size": 0,
5   "aggs": {
6     "job_agg": {
7       "terms": {
8         "field": "job.keyword"
9       }
10    },
11    "max_salary": {
12      "max": {
13        "field": "salary"
14      }
15    }
16  }
17 }
18 # msearch 拆分多个语句的聚合实现
19 GET _msearch
20 {"index":"employees"}
21 {"size":0,"aggs":{"job_agg":{"terms":{"field": "job.keyword"}}}}
22 {"index":"employees"}
23 {"size":0,"aggs":{"max_salary":{"max":{"field": "salary"}}}}
```