

主讲老师: Fox

- 1 文档: 3. Elasticsearch搜索技术深入与聚合查...
- 2 链接: <http://note.youdao.com/noteshare?id=44848fef10c56f33bdb7a85f7ca04376&sub=89B65B40D9584654978759587D3ADEB9>

相关性和相关性算分

相关性 (Relevance)

什么是TF-IDF

BM25

通过Explain API查看TF-IDF

Boosting

布尔查询bool Query

如何解决结构化查询“包含而不是相等”的问题

利用bool嵌套实现should not逻辑

Boosting Query

控制字段的Boosting

单字符串多字段查询

三种场景

最佳字段查询Dis Max Query

Multi Match Query

ElasticSearch聚合操作

聚合的分类

Metric Aggregation

Bucket Aggregation

Pipeline Aggregation

聚合的作用范围

排序

ES聚合分析不精准原因分析

Elasticsearch 聚合性能优化

启用 eager global ordinals 提升高基数聚合性能

插入数据时对索引进行预排序

使用节点查询缓存

使用分片请求缓存

拆分聚合，使聚合并行化

相关性和相关性算分

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

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

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。**ES 5之前，默认的相关性算分采用TF-IDF，现在采用BM 25。

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

什么是TF-IDF

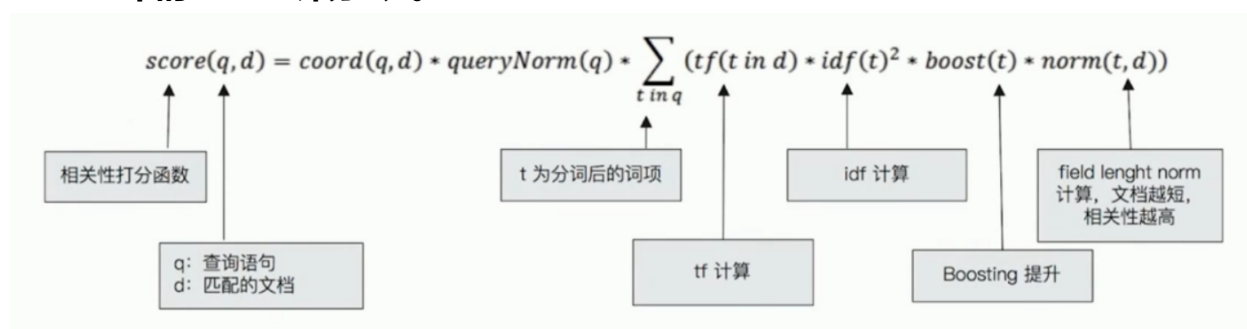
TF-IDF (term frequency-inverse document frequency) 是一种用于信息检索与数据挖掘的常用加权技术。

- 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(全部文档数/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

- 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。

- **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。

- **字段长度归一值 (field-length norm)**

字段的长度是多少？字段越短，字段的权重越高。检索词出现在一个内容短 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)、逆向文档频率 (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

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

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 formul
a"}
```

```

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 }

```

Boosting

Boosting是控制相关度的一种手段。

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

返回匹配positive查询的文档并降低匹配negative查询的文档相似度分。这样就可以在不排除某些文档的前提下对文档进行查询,搜索结果中存在只不过相似度分数相比正常匹配的要低:

```

1 GET /test_score/_search
2 {
3   "query": {
4     "boosting": {
5       "positive": {
6         "term": {
7           "content": "elasticsearch"
8         }
9       },
10      "negative": {
11        "term": {
12          "content": "like"
13        }
14      },
15      "negative_boost": 0.2

```

```
16  }  
17  }  
18  }
```

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。

布尔查询bool Query

一个bool查询,是一个或者多个查询子句的组合，总共包括4种子句，其中2种会影响算分，2种不影响算分。

- must: 相当于&&，必须匹配，贡献算分
- should: 相当于||，选择性匹配，贡献算分
- must_not: 相当于!，必须不能匹配，不贡献算分
- filter: 必须匹配，不贡献算法

在Elasticsearch中，有Query和 Filter两种不同的Context

- Query Context: 相关性算分
- Filter Context: 不需要算分,可以利用Cache，获得更好的性能

相关性并不只是全文本检索的专利，也适用于yes | no 的子句，匹配的子句越多，相关性评分

越高。如果多条查询子句被合并为一条复合查询语句，比如 bool查询，则每个查询子句计算得出的评分会被合并到总的相关性评分中。

bool查询语法

- 子查询可以任意顺序出现
- 可以嵌套多个查询
- 如果你的bool查询中，没有must条件,should中必须至少满足一条查询

```
1 GET /es_db/_search  
2 {  
3   "query": {  
4     "bool": {  
5       "must": {  
6         "match": {  
7           "remark": "java developer"  
8         }  
9       },  
10      "filter": {  
11        "term": {
```

```
12  "sex": "1"
13  }
14  },
15  "must_not": {
16    "range": {
17      "age": {
18        "gte": 30
19      }
20    }
21  },
22  "should": [
23    {
24      "term": {
25        "address.keyword": {
26          "value": "广州天河公园"
27        }
28      }
29    },
30    {
31      "term": {
32        "address.keyword": {
33          "value": "广州白云山公园"
34        }
35      }
36    }
37  ],
38  "minimum_should_match": 1
39  }
40  }
41  }
```

如何解决结构化查询“包含而不是相等”的问题

测试数据

```
1  POST /employee/_bulk
2  {"index":{"_id":1}}
3  {"name":"小明","interest":["跑步","篮球"]}
4  {"index":{"_id":2}}
5  {"name":"小红","interest":["跑步"]}
6  {"index":{"_id":3}}
7  {"name":"小丽","interest":["跳舞","唱歌","跑步"]}
```

```

8
9 POST /employee/_search
10 {
11   "query": {
12     "term": {
13       "interest.keyword": {
14         "value": "跑步"
15       }
16     }
17   }
18 }

```

解决方案：增加count字段，使用bool查询解决

- 从业务角度，按需改进Elasticsearch数据模型

```

1 POST /employee/_bulk
2 {"index":{"_id":1}}
3 {"name":"小明","interest":["跑步","篮球"],"interest_count":2}
4 {"index":{"_id":2}}
5 {"name":"小红","interest":["跑步"],"interest_count":1}
6 {"index":{"_id":3}}
7 {"name":"小丽","interest":["跳舞","唱歌","跑步"],"interest_count":3}

```

- 使用bool查询

```

1 # must 算分
2 POST /employee/_search
3 {
4   "query": {
5     "bool": {
6       "must": [
7         {
8           "term": {
9             "interest.keyword": {
10              "value": "跑步"
11            }
12          }
13        },
14        {
15          "term": {
16            "interest_count": {
17              "value": 1
18            }
19          }
20        }
21      ]
22    }
23  }
24 }

```



```

19  }
20  }
21  ]
22  }
23  }
24  }
25  # filter不算分
26  POST /employee/_search
27  {
28    "query": {
29      "bool": {
30        "filter": [
31          {
32            "term": {
33              "interest.keyword": {
34                "value": "跑步"
35              }
36            }
37          },
38          {
39            "term": {
40              "interest_count": {
41                "value": 1
42              }
43            }
44          }
45        ]
46      }
47    }
48  }

```

利用bool嵌套实现should not逻辑

```

1  GET /es_db/_search
2  {
3    "query": {
4      "bool": {
5        "must": {
6          "match": {
7            "remark": "java developer"
8          }

```

```

9   },
10  "should": [
11    {
12      "bool": {
13        "must_not": [
14          {
15            "term": {
16              "sex": 1
17            }
18          }
19        ]
20      }
21    }
22  ],
23  "minimum_should_match": 1
24 }
25 }
26 }

```

Boosting Query

思考：如何控制查询的相关性算分？

控制字段的Boosting

Boosting是控制相关的一种手段。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

```

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
23
24
25
```

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

利用negative_boost降低相关性

- 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

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.

```
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 }
```

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

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 PUT /blogs/_doc/1
```

```

2 {
3   "title": "Quick brown rabbits",
4   "body": "Brown rabbits are commonly seen."
5 }
6
7 PUT /blogs/_doc/2
8 {
9   "title": "Keeping pets healthy",
10  "body": "My quick brown fox eats rabbits on a regular basis."
11 }
12
13 POST /blogs/_search
14 {
15   "query": {
16     "bool": {
17       "should": [
18         { "match": { "title": "Brown fox" }},
19         { "match": { "body": "Brown fox" }}
20       ]
21     }
22   }
23 }
24

```

思考：查询结果不符合预期，为什么？

<pre> 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" }}] } } } </pre>	<pre> 13 "relation": "eq" 14 }, 15 "max_score" : 0.90425634, 16 "hits" : [17 { 18 "_index" : "blogs", 19 "_type" : "_doc", 20 "_id" : "1", 21 "_score" : 0.90425634, 22 "_source" : { 23 "title" : "Quick brown rabbits", 24 "body" : "Brown rabbits are commonly seen." 25 } 26 }, 27 { 28 "_index" : "blogs", 29 "_type" : "_doc", 30 "_id" : "2", 31 "_score" : 0.77041256, 32 "_source" : { 33 "title" : "Keeping pets healthy", 34 "body" : "My quick brown fox eats rabbits on a regular basis." 35 } 36 } 37] 38 } </pre>
---	--

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，不应该讲分数简单叠加，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": 0.2
23 }
24 }
25 }
```

Multi Match Query

最佳字段(Best Fields)搜索

Best Fields是默认类型，可以不用指定

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

使用多数字段 (Most Fields) 搜索

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 }
```

跨字段 (Cross Field) 搜索

```
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
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
46 GET /address/_search
47 {
48   "query": {
49     "multi_match": {
50       "query": "湖南常德",
51       "type": "most_fields",
52       "fields": ["province","city"]
53     }
54   }
55 }
```

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?

语法:

```
1 "aggs" : { #和query同级的关键词
2   "<aggregation_name>" : { #自定义的聚合名字
3     "<aggregation_type>" : { #聚合的定义: 不同的type+body
4       <aggregation_body>
5     }
6   [, "meta" : { [<meta_data_body>] } ]?
7   [, "aggregations" : { [<sub_aggregation>]+ } ]? #子聚合查询
8 }
9 [, "<aggregation_name_2>" : { ... } ]* #可以包含多个同级的聚合查询
10 }
```

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 SELECT 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}

```

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  }
```

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  }
12 }
```

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   }
```

```
10 }
11 }
```

```
{
  "error" : {
    "root_cause" : [
      {
        "type" : "illegal_argument_exception",
        "reason" : "Text fields are not optimised for operations that require per-document field data like aggregations and sorting, so these operations are disabled by default. Please use a keyword field instead. Alternatively, set fielddata=true on [job] in order to load field data by uninverting the inverted index. Note that this can use significant memory."
      }
    ],
    "type" : "search_phase_execution_exception",
    "reason" : "all shards failed",
    "phase" : "query",
    "grouped" : true,
    "failed_shards" : [
```

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 }
```

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  }
```

```
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 }
46 }

```

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   "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
```

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   }
44 }
```

```
42  "all_jobs": {
43    "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

```

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   },

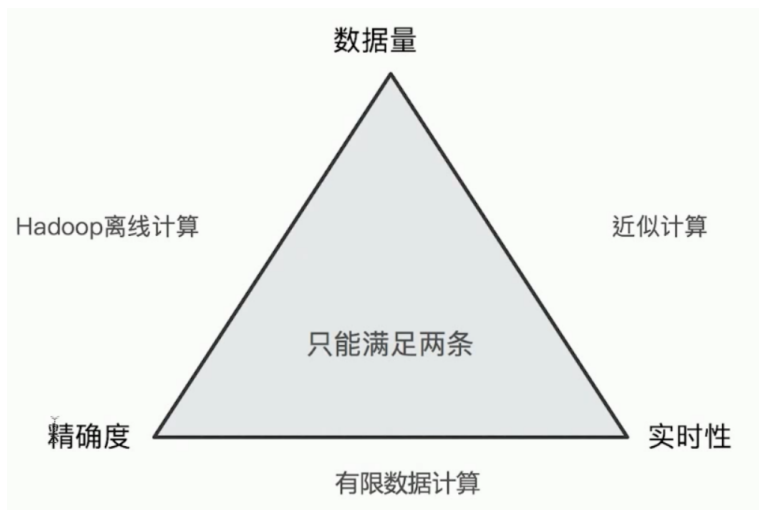
```

```
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
```

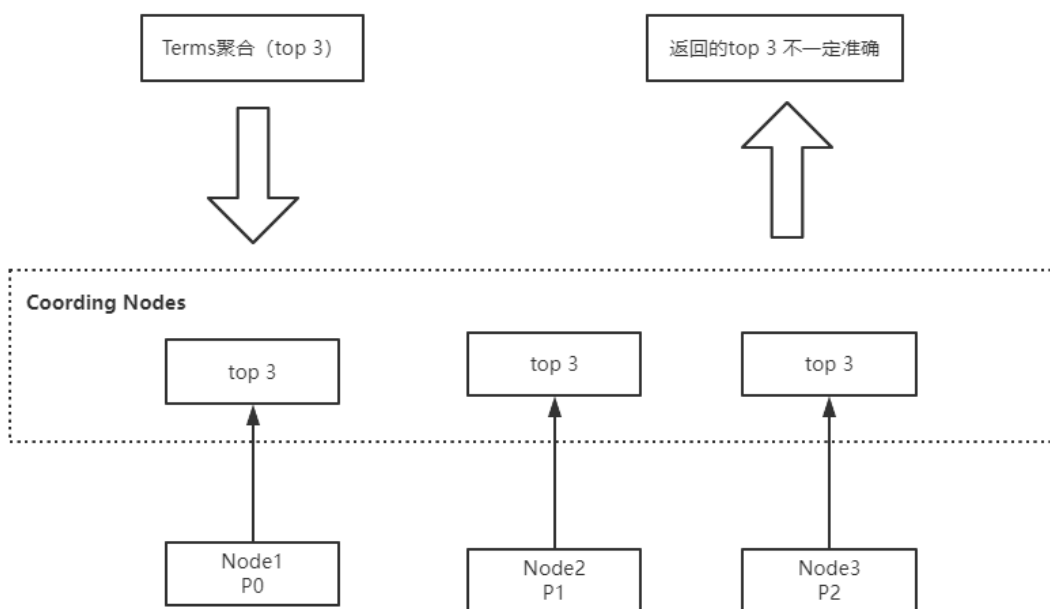
```
54
55 #排序 order
56 #count and key
57 POST employees/_search
58 {
59   "size": 0,
60   "aggs": {
61     "jobs": {
62       "terms": {
63         "field": "job.keyword",
64         "order": [ {
65           "stats_salary.min": "desc"
66         } ]
67       }
68     },
69     "aggs": {
70       "stats_salary":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

ES聚合分析不精准原因分析

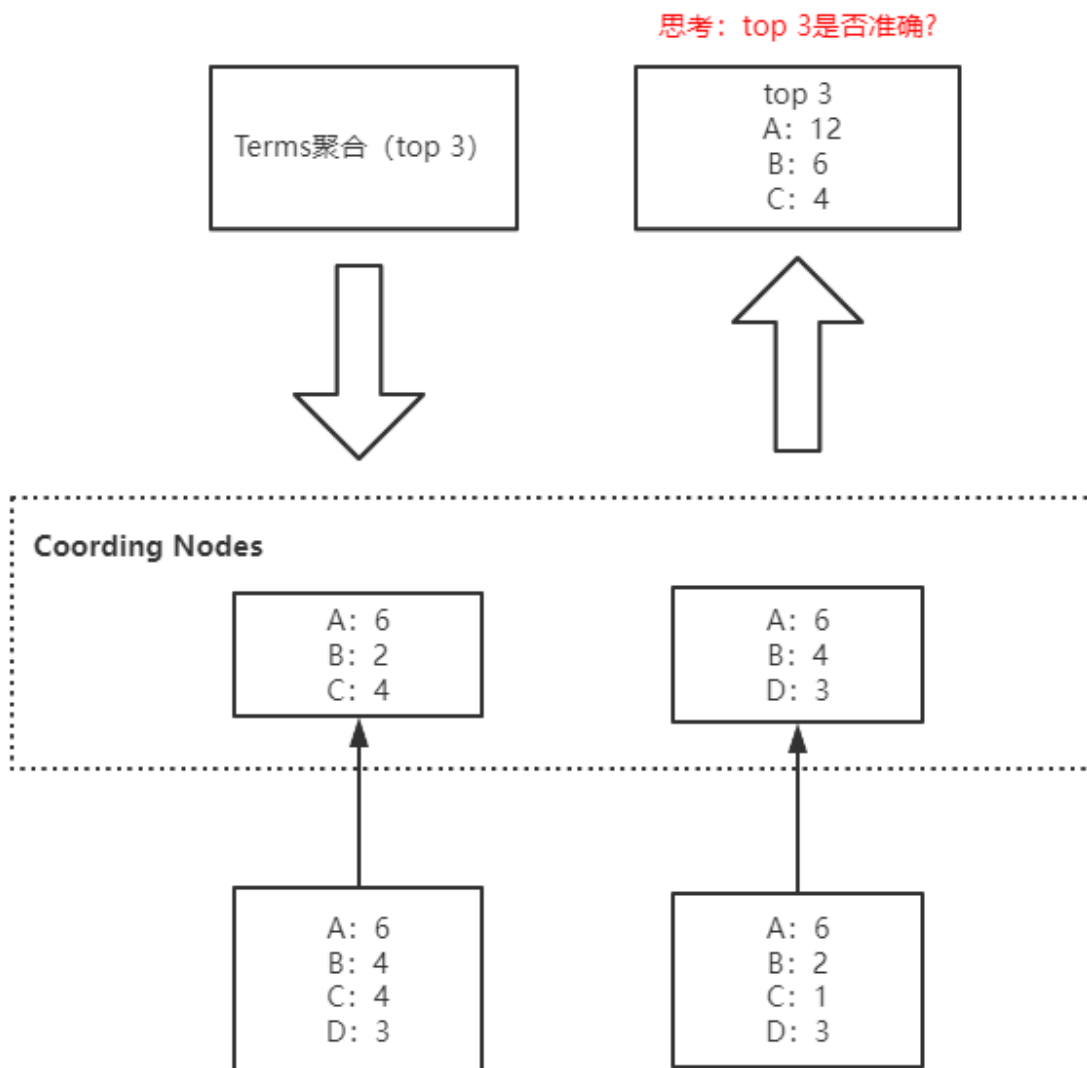
ElasticSearch在对海量数据进行聚合分析的时候会损失搜索的精准度来满足实时性的需求。



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： $\text{size} \times 1.5 + 10$ 。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。

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  }
```

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   }
6 }
```

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  }
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"}}}}
```