

主讲老师: Fox

有道笔记地址: <https://note.youdao.com/s/3lIHx6vE>

1. ES高级查询Query DSL

ES中提供了一种强大的检索数据方式,这种检索方式称之为Query DSL (Domain Specified Language 领域专用语言) , Query DSL是利用Rest API传递JSON格式的请求体(RequestBody)数据与ES进行交互, 这种方式的丰富查询语法让ES检索变得更强大, 更简洁。

<https://www.elastic.co/guide/en/elasticsearch/reference/7.17/query-dsl.html>

语法:

```
1 GET /es_db/_doc/_search {json请求体数据}
2 可以简化为下面写法
3 GET /es_db/_search {json请求体数据}
```

示例

```
1 #无条件查询, 默认返回10条数据
2 GET /es_db/_search
3 {
4     "query":{
5         "match_all":{}
6     }
7 }
```

```

1  # [types removal] specifying types in search requests is deprecated
2  {
3    "took" : 1,
4    "timed_out" : false,
5    "_shards" : {
6      "total" : 1,
7      "successful" : 1,
8      "skipped" : 0,
9      "failed" : 0
10   },
11   "hits" : {
12     "total" : {
13       "value" : 13,
14       "relation" : "eq"
15     },
16     "max_score" : 1.0,
17     "hits" : [
18       {
19         "_index" : "es_db",
20         "_type" : "_doc",
21         "_id" : "2",
22         "_score" : 1.0,
23         "_source" : {
24           "name" : "李四",
25           "sex" : 1,
26           "age" : 28,
27           "address" : "广州荔湾大厦",
28           "remark" : "java assistant"
29         }
30       }
31     ]
32   }
33 }

```

took:花费的时间

total.value: 符合条件的总文档

hits: 结果集, 默认前10个文档

_index:索引名
_id: 文档的id
_score: 相关度评分
source:文档原生信息

示例数据

```

1  #指定ik分词器
2  PUT /es_db
3  {
4    "settings" : {
5      "index" : {
6        "analysis.analyzer.default.type": "ik_max_word"
7      }
8    }
9  }
10
11 # 创建文档,指定id
12 PUT /es_db/_doc/1
13 {
14   "name": "张三",

```

```
15  "sex": 1,
16  "age": 25,
17  "address": "广州天河公园",
18  "remark": "java developer"
19  }
20  PUT /es_db/_doc/2
21  {
22  "name": "李四",
23  "sex": 1,
24  "age": 28,
25  "address": "广州荔湾大厦",
26  "remark": "java assistant"
27  }
28
29  PUT /es_db/_doc/3
30  {
31  "name": "王五",
32  "sex": 0,
33  "age": 26,
34  "address": "广州白云山公园",
35  "remark": "php developer"
36  }
37
38  PUT /es_db/_doc/4
39  {
40  "name": "赵六",
41  "sex": 0,
42  "age": 22,
43  "address": "长沙橘子洲",
44  "remark": "python assistant"
45  }
46
47  PUT /es_db/_doc/5
48  {
49  "name": "张龙",
50  "sex": 0,
51  "age": 19,
52  "address": "长沙麓谷企业广场",
53  "remark": "java architect assistant"
54  }
```

```
55
56 PUT /es_db/_doc/6
57 {
58   "name": "赵虎",
59   "sex": 1,
60   "age": 32,
61   "address": "长沙麓谷兴工国际产业园",
62   "remark": "java architect"
63 }
64
65 PUT /es_db/_doc/7
66 {
67   "name": "李虎",
68   "sex": 1,
69   "age": 32,
70   "address": "广州番禺节能科技园",
71   "remark": "java architect"
72 }
73
74 PUT /es_db/_doc/8
75 {
76   "name": "张星",
77   "sex": 1,
78   "age": 32,
79   "address": "武汉东湖高新区未来智汇城",
80   "remark": "golang developer"
81 }
82
```

1.1 match_all

使用match_all，匹配所有文档，默认只会返回10条数据。

原因：_search查询默认采用的是分页查询，每页记录数size的默认值为10。如果想显示更多数据，指定size

```
1 GET /es_db/_search
2 等同于
```

```
3 GET /es_db/_search
4 {
5   "query":{
6     "match_all":{}
7   }
8 }
```

返回源数据_source

_source 关键字: 是一个数组,在数组中用来指定展示那些字段

```
1 # 返回指定字段
2 GET /es_db/_search
3 {
4   "query": {
5     "match_all": {}
6   },
7   "_source": ["name", "address"]
8 }
9
10 #在查询中过滤
11 #不查看源数据, 仅查看元字段
12 {
13   "_source": false,
14   "query": {
15     ...
16   }
17 }
18
19 #只看以obj.开头的字段
20 {
21   "_source": "obj.*",
22   "query": {
23     ...
24   }
25 }
```

返回指定条数size

size 关键字: 指定查询结果中返回指定条数。 默认返回值10条

```
1 GET /es_db/_search
2 {
3   "query": {
4     "match_all": {}
5   },
6   "size": 100
7 }
8
```

分页查询from&size

size: 显示应该返回的结果数量, 默认是 10

from: 显示应该跳过的初始结果数量, 默认是 0

from 关键字用来指定起始返回位置, 和size关键字连用可实现分页效果

```
1 GET /es_db/_search
2 {
3   "query": {
4     "match_all": {}
5   },
6   "from": 0,
7   "size": 5
8 }
```

指定字段排序sort

注意: 会让得分失效

```
1 GET /es_db/_search
```

```
2  {
3    "query": {
4      "match_all": {}
5    },
6    "sort": [
7      {
8        "age": "desc"
9      }
10   ]
11 }
12
13 #排序, 分页
14 GET /es_db/_search
15 {
16   "query": {
17     "match_all": {}
18   },
19   "sort": [
20     {
21       "age": "desc"
22     }
23   ],
24   "from": 10,
25   "size": 5
26 }
```

1.2 术语级别查询

术语级别查询 (Term-Level Queries) 指的是搜索内容不经过文本分析直接用于文本匹配, 这个过程类似于数据库的SQL查询, 搜索的对象大多是索引的非text类型字段。Elasticsearch 中的一些术语级别查询示例包括 term、terms 和 range 查询。

Term query术语查询

术语查询直接返回包含搜索内容的文档, 常用来查询索引中某个类型为keyword的文本字段, 类似于SQL的 “=” 查询, 使用十分普遍。

注意: 最好不要在term查询的字段中使用text字段, 因为text字段会被分词, 这样做既没有意义, 还很有可能什么也查不到。

```
1 # 对bool, 日期, 数字, 结构化的文本可以利用term做精确匹配
2 # term 精确匹配
3 GET /es_db/_search
4 {
5   "query": {
6     "term": {
7       "age": {
8         "value": 28
9       }
10    }
11  }
12 }
13
14 # 思考: 查询广州白云是否有数据, 为什么?
15 GET /es_db/_search
16 {
17   "query":{
18     "term": {
19       "address": {
20         "value": "广州白云"
21       }
22     }
23   }
24 }
25
26 # 采用term精确查询, 查询字段映射类型为keyword
27 GET /es_db/_search
28 {
29   "query":{
30     "term": {
31       "address.keyword": {
32         "value": "广州白云山公园"
33       }
34     }
35   }
36 }
37
38
```


在ES中，Term查询，对输入不做分词。会将输入作为一个整体，在倒排索引中查找准确的词项，并且使用相关度算分公式为每个包含该词项的文档进行相关度算分。

可以通过 Constant Score 将查询转换成一个 Filtering，避免算分，并利用缓存，提高性能。

- 将Query 转成 Filter，忽略TF-IDF计算，避免相关性算分的开销
- Filter可以有效利用缓存

```
1 GET /es_db/_search
2 {
3   "query": {
4     "constant_score": {
5       "filter": {
6         "term": {
7           "address.keyword": "广州白云山公园"
8         }
9       }
10    }
11  }
12 }
```

term处理多值字段时，term查询是包含，不是等于。

```
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":["跳舞","画画"]}
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 }
```

```
17     }  
18 }
```

Terms Query多术语查询

Terms query用于在指定字段上匹配多个词项（terms）。它会精确匹配指定字段中包含的任何一个词项。

```
1 POST /es_db/_search  
2 {  
3   "query": {  
4     "terms": {  
5       "remark.keyword": ["java assistant", "java architect"]  
6     }  
7   }  
8 }  
9
```

exists query

在Elasticsearch中可以使用exists进行查询，以判断文档中是否存在对应的字段。

```
1 #查询索引库中存在remarks字段的文档数据  
2 GET /es_db/_search  
3 {  
4   "query": {  
5     "exists": {  
6       {  
7         "field": "remark"  
8       }  
9     }  
10  }
```

ids query

ids 关键字：值为数组类型,用来根据一组id获取多个对应的文档

```
1 GET /es_db/_search
2 {
3   "query": {
4     "ids": {
5       "values": [1,2]
6     }
7   }
8 }
```

range query范围查询

- range: 范围关键字
- gte 大于等于
- lte 小于等于
- gt 大于
- lt 小于
- now 当前时间

```
1 POST /es_db/_search
2 {
3   "query": {
4     "range": {
5       "age": {
6         "gte": 25,
7         "lte": 28
8       }
9     }
10  }
11 }
12
13 #日期范围比较
14 DELETE /product
15 POST /product/_bulk
16 {"index":{"_id":1}}
17 {"price":100,"date":"2021-01-01","productId":"XHDK-1293"}
18 {"index":{"_id":2}}
```

```
19  {"price":200,"date":"2022-01-01","productId":"KDKE-5421"}
20
21  GET /product/_mapping
22
23  GET /product/_search
24  {
25    "query": {
26      "range": {
27        "date": {
28          "gte": "now-2y"
29        }
30      }
31    }
32  }
33
34
```

prefix query前缀查询

它会对分词后的term进行前缀搜索。

- 它不会分析要搜索字符串，传入的前缀就是想要查找的前缀
- 默认状态下，前缀查询不做相关度分数计算，它只是将所有匹配的文档返回，然后赋予所有相关分数值为1。它的行为更像是一个过滤器而不是查询。两者实际的区别就是过滤器是可以被缓存的，而前缀查询不行。

prefix的原理：需要遍历所有倒排索引，并比较每个term是否以所指定的前缀开头。

```
1  GET /es_db/_search
2  {
3    "query": {
4      "prefix": {
5        "address": {
6          "value": "广州"
7        }
8      }
9    }
10 }
```

wildcard query通配符查询

通配符查询：工作原理和prefix相同，只不过它不是只比较开头，它能支持更为复杂的匹配模式。

```
1 GET /es_db/_search
2 {
3   "query": {
4     "wildcard": {
5       "address": {
6         "value": "*白*"
7       }
8     }
9   }
10 }
```

fuzzy query模糊查询

在实际的搜索中，我们有时候会打错字，从而导致搜索不到。在Elasticsearch中，我们可以使用fuzziness属性来进行模糊查询，从而达到搜索有错别字的情形。

fuzzy 查询会用到两个很重要的参数，fuzziness, prefix_length

- fuzziness：表示输入的关键词通过几次操作可以转变成成为ES库里面的对应field的字段
 - 操作是指：新增一个字符，删除一个字符，修改一个字符，每次操作可以记做编辑距离为1；
 - 如中文集团到中威集团编辑距离就是1，只需要修改一个字符；如果fuzziness值在这里设置成2，会把编辑距离为2的东东集团也查出来。
 - 该参数默认值为0，即不开启模糊查询；fuzzy 模糊查询 最大模糊错误必须在0-2之间
- prefix_length：表示限制输入关键字和ES对应查询field的内容开头的第n个字符必须完全匹配，不允许错别字匹配；
 - 如这里等于1，则表示开头的字必须匹配，不匹配则不返回；
 - 默认值也是0；
 - 加大prefix_length的值可以提高效率和准确率。

```
1 GET /es_db/_search
2 {
3   "query": {
4     "fuzzy": {
5       "address": {
6         "value": "白运山",
7         "fuzziness": 1
8       }
9     }
10  }
11 }
12
```

1.3 全文检索

全文检索查询（Full Text Queries）和术语级别查询（Term-Level Queries）是 Elasticsearch 中搜索和检索数据的两种不同方法。

全文检索查询旨在基于相关性搜索和匹配文本数据。这些查询会对输入的文本进行分析，将其拆分为词项（单个单词），并执行诸如分词、词干处理和标准化等操作。Elasticsearch 中的一些全文检索查询示例包括 match、match_phrase 和 multi_match 查询。

全文检索的关键特点：

- **对输入的文本进行分析，并根据分析后的词项进行搜索和匹配。**全文检索查询会对输入的文本进行分析，将其拆分为词项，并基于这些词项进行搜索和匹配操作。
- 以相关性为基础进行搜索和匹配。全文检索查询使用相关性算法来确定文档与查询的匹配程度，并按照相关性进行排序。相关性可以基于词项的频率、权重和其他因素来计算。
- 全文检索查询适用于包含自由文本数据的字段，例如文档的内容、文章的正文或产品描述等。

match query匹配查询

match在匹配时会对所查找的关键词进行分词，然后按分词匹配查找。

match支持以下参数：

- query : 指定匹配的值
- operator : 匹配条件类型
 - and : 条件分词后都要匹配
 - or : 条件分词后有一个匹配即可(默认)

- `minnum_should_match`: 最低匹配度, 即条件在倒排索引中最低的匹配度

```
1 #match 分词后or的效果
2 GET /es_db/_search
3 {
4   "query": {
5     "match": {
6       "address": "广州白云山公园"
7     }
8   }
9 }
10
11 # 分词后 and的效果
12 GET /es_db/_search
13 {
14   "query": {
15     "match": {
16       "address": {
17         "query": "广州白云山公园",
18         "operator": "and"
19       }
20     }
21   }
22 }
23
```

在match中的应用: 当operator参数设置为or时, `minnum_should_match`参数用来控制匹配的分词的最少数量。

```
1 # 最少匹配广州, 公园两个词
2 GET /es_db/_search
3 {
4   "query": {
5     "match": {
6       "address": {
7         "query": "广州公园",
```

```
8         "minimum_should_match": 2
9     }
10 }
11 }
12 }
```

对于match查询，其底层逻辑的概述：

1. 分词：首先，输入的查询文本会被分词器进行分词。分词器会将文本拆分成一个个词项（terms），如单词、短语或特定字符。分词器通常根据特定的语言规则和配置进行操作。
2. 倒排索引：ES使用倒排索引来加速搜索过程。倒排索引是一种数据结构，它将词项映射到包含这些词项的文档。每个词项都有一个对应的倒排列表，其中包含了包含该词项的所有文档的引用。
3. 匹配计算：一旦查询被分词，ES将根据查询的类型和参数计算文档与查询的匹配度。对于match查询，ES将比较查询的词项与倒排索引中的词项，并计算文档的相关性得分。相关性得分衡量了文档与查询的匹配程度。
4. 结果返回：根据相关性得分，ES将返回最匹配的文档作为搜索结果。搜索结果通常按照相关性得分进行排序，以便最相关的文档排在前面。

multi_match query 多字段查询

可以根据字段类型，决定是否使用分词查询，得分最高的在前面

```
1 GET /es_db/_search
2 {
3   "query": {
4     "multi_match": {
5       "query": "长沙张龙",
6       "fields": [
7         "address",
8         "name"
9       ]
10    }
11  }
12 }
```

注意：字段类型分词,将查询条件分词之后进行查询，如果该字段不分词就会将查询条件作为整体进行查询。

match_phrase query短语查询

短语搜索(match phrase)会对搜索文本进行文本分析，然后到索引中寻找搜索的每个分词并**要求分词相邻**，你可以通过调整slop参数设置分词出现的最大间隔距离。**match_phrase** 会将检索关键词分词。

```
1 GET /es_db/_search
2 {
3   "query": {
4     "match_phrase": {
5       "address": "广州白云山"
6     }
7   }
8 }
9 GET /es_db/_search
10 {
11   "query": {
12     "match_phrase": {
13       "address": "广州白云"
14     }
15   }
16 }
```

思考：为什么查询广州白云山有数据，广州白云没有数据？

```
3 ^ | }
3 ^ | }
3 ^ | }
1 ^ }
2
3 # match phrase 短语匹配
4
5 GET /es_db/_search
6 {
7   "query": {
8     "match_phrase": {
9       "address": "广州白云山"
10     }
11   }
12 }
13
14 GET /es_db/_search
15 {
16   "query": {
17     "match_phrase": {
18       "address": "广州白云"
19     }
20   }
21 }
22
23 GET /es_db/_search
24 {
25   "query": {
26     "match_phrase": {
27       "address": "广州云山",
28       "slop": 2
29     }
30   }
31 }
```

```
15 ^ },
16   "max_score" : 4.949977,
17   "hits" : [
18     {
19       "_index" : "es_db",
20       "_type" : "_doc",
21       "_id" : "3",
22       "_score" : 4.949977,
23       "_source" : {
24         "name" : "王五",
25         "sex" : 0,
26         "age" : 26,
27         "address" : "广州白云山公园",
28         "remark" : "php developer"
29       }
30     }
31   ]
32 }
33
34 # GET /es_db/_search
35 {
36   "took" : 1,
37   "timed_out" : false,
38   "_shards" : { },
39   "hits" : {
40     "total" : {
41       "value" : 0,
42       "relation" : "eq"
43     },
44     "max_score" : null,
45     "hits" : [ ]
46   }
47 }
```

思考：为什么广州白云山有数据，查广州白云没数据？

分析原因：

先查看广州白云山公园分词结果，可以知道广州和白云不是相邻的词条，中间会隔一个白云山，而match_phrase匹配的是相邻的词条，所以查询广州白云山有结果，但查询广州白云没有结果。

```
1 POST _analyze
2 {
3     "analyzer":"ik_max_word",
4     "text":"广州白云山"
5 }
6 #结果
7 {
8     "tokens" : [
9         {
10             "token" : "广州",
11             "start_offset" : 0,
12             "end_offset" : 2,
13             "type" : "CN_WORD",
14             "position" : 0
15         },
16         {
17             "token" : "白云山",
18             "start_offset" : 2,
19             "end_offset" : 5,
20             "type" : "CN_WORD",
21             "position" : 1
22         },
23         {
24             "token" : "白云",
25             "start_offset" : 2,
26             "end_offset" : 4,
27             "type" : "CN_WORD",
28             "position" : 2
29         },
30         {
31             "token" : "云山",
32             "start_offset" : 3,
33             "end_offset" : 5,
34             "type" : "CN_WORD",
35             "position" : 3
36         }
37     ]
38 }
```

```
37     ]
38 }
```

如何解决词条间隔的问题？可以借助slop参数，slop参数告诉match_phrase查询词条能够相隔多远时仍然将文档视为匹配。

```
1  #广州云山分词后相隔为2，可以匹配到结果
2  GET /es_db/_search
3  {
4    "query": {
5      "match_phrase": {
6        "address": {
7          "query": "广州云山",
8          "slop": 2
9        }
10     }
11  }
12 }
13
```

query_string query

允许我们在单个查询字符串中指定AND | OR | NOT条件，同时也和 multi_match query 一样，支持多字段搜索。和match类似，但是match需要指定字段名，query_string是在所有字段中搜索，范围更广泛。

注意: 查询字段分词就将查询条件分词查询，查询字段不分词将查询条件不分词查询

- 未指定字段查询

```
1  # AND 要求大写
2  GET /es_db/_search
3  {
4    "query": {
5      "query_string": {
6        "query": "赵六 AND 橘子洲"
7      }
8    }
9  }
```

```
8   }  
9 }
```

- 指定单个字段查询

```
1 #Query String  
2 GET /es_db/_search  
3 {  
4   "query": {  
5     "query_string": {  
6       "default_field": "address",  
7       "query": "白云山 OR 橘子洲"  
8     }  
9   }  
10 }
```

- 指定多个字段查询

```
1 GET /es_db/_search  
2 {  
3   "query": {  
4     "query_string": {  
5       "fields": ["name", "address"],  
6       "query": "张三 OR (广州 AND 王五)"  
7     }  
8   }  
9 }  
10
```

simple_query_string

类似Query String，但是会忽略错误的语法,同时只支持部分查询语法，不支持AND OR NOT，会当作字符串处理。支持部分逻辑：

- + 替代AND

- | 替代OR
- - 替代NOT

```
1 #simple_query_string 默认的operator是OR
2 GET /es_db/_search
3 {
4   "query": {
5     "simple_query_string": {
6       "fields": ["name","address"],
7       "query": "广州公园",
8       "default_operator": "AND"
9     }
10  }
11 }
12
13 GET /es_db/_search
14 {
15   "query": {
16     "simple_query_string": {
17       "fields": ["name","address"],
18       "query": "广州 + 公园"
19     }
20  }
21 }
```

1.4 bool query布尔查询

布尔查询可以按照布尔逻辑条件组织多条查询语句，只有符合整个布尔条件的文档才会被搜索出来。在布尔条件中，可以包含两种不同的上下文。

1. **搜索上下文(query context)**: 使用搜索上下文时，Elasticsearch需要计算每个文档与搜索条件的相关度得分，这个得分的计算需使用一套复杂的计算公式，**有一定的性能开销，带文本分析的全文检索的查询语句很适合放在搜索上下文中。**
2. **过滤上下文(filter context)**: 使用过滤上下文时，Elasticsearch只需要判断搜索条件跟文档数据是否匹配，例如使用Term query判断一个值是否跟搜索内容一致，使用Range query判断某数据是否位于某个区间等。**过滤上**

下文的查询不需要进行相关度得分计算，还可以使用缓存加快响应速度，很多术语级查询语句都适合放在过滤上下文中。

布尔查询一共支持4种组合类型:

类型	说明
must	可包含多个查询条件，每个条件均满足的文档才能被搜索到，每次查询需要计算相关度得分，属于搜索上下文
should	可包含多个查询条件，不存在must和fiter条件时，至少要满足多个查询条件中的一个，文档才能被搜索到，否则需满足的条件数量不受限制,匹配到的查询越多相关度越高，也属于搜索上下文
filter	可包含多个过滤条件，每个条件均满足的文档才能被搜索到，每个过滤条件不计算相关度得分，结果在一定条件下会被缓存，属于过滤上下文
must_not	可包含多个过滤条件，每个条件均不满足的文档才能被搜索到，每个过滤条件不计算相关度得分，结果在一定条件下会被缓存，属于过滤上下文

示例

```
1 PUT /books
2 {
3   "settings": {
4     "number_of_replicas": 1,
5     "number_of_shards": 1
6   },
7   "mappings": {
8     "properties": {
9       "id": {
10        "type": "long"
11      },
12      "title": {
13        "type": "text",
14        "analyzer": "ik_max_word"
15      },
16      "language": {
17        "type": "keyword"
```

```
18     },
19     "author": {
20         "type": "keyword"
21     },
22     "price": {
23         "type": "double"
24     },
25     "publish_time": {
26         "type": "date",
27         "format": "yyy-MM-dd"
28     },
29     "description": {
30         "type": "text",
31         "analyzer": "ik_max_word"
32     }
33 }
34 }
35 }
36
37 POST /_bulk
38 {"index":{"_index":"books","_id":"1"}}
39 {"id":"1","title":"Java编程思想","language":"java","author":"Bruce Eckel",
40  "price":70.20,"publish_time":"2007-10-01","description":"Java学习必读经典，殿堂级著作！
41  赢得了全球程序员的广泛赞誉。"}
42 {"index":{"_index":"books","_id":"2"}}
43 {"id":"2","title":"Java程序性能优化","language":"java","author":"葛一
44  鸣","price":46.5,"publish_time":"2012-08-01","description":"让你的Java程序更快、更稳定。
45  深入剖析软件设计层面、代码层面、JVM虚拟机层面的优化方法"}
46 {"index":{"_index":"books","_id":"3"}}
47 {"id":"3","title":"Python科学计算","language":"python","author":"张若
48  愚","price":81.4,"publish_time":"2016-05-01","description":"零基础学python，光盘中作者独
49  家整合开发winPython运行环境，涵盖了Python各个扩展库"}
50 {"index":{"_index":"books","_id":"4"}}
51 {"id":"4","title":"Python基础教程","language":"python","author":"Helant",
52  "price":54.50,"publish_time":"2014-03-01","description":"经典的Python入门教程，层次鲜
53  明，结构严谨，内容翔实"}
54 {"index":{"_index":"books","_id":"5"}}
55 {"id":"5","title":"JavaScript高级程序设计","language":"javascript","author":"Nicholas
56  C. Zakas","price":66.4,"publish_time":"2012-10-01","description":"JavaScript技术经典名
57  著"}
58
59 GET /books/_search
```

```
51 {
52   "query": {
53     "bool": {
54       "must": [
55         {
56           "match": {
57             "title": "java编程"
58           }
59         }, {
60           "match": {
61             "description": "性能优化"
62           }
63         }
64       ]
65     }
66   }
67 }
```

```
68
69 GET /books/_search
```

```
70 {
71   "query": {
72     "bool": {
73       "should": [
74         {
75           "match": {
76             "title": "java编程"
77           }
78         }, {
79           "match": {
80             "description": "性能优化"
81           }
82         }
83       ],
84       "minimum_should_match": 1
85     }
86   }
87 }
```

```
88 GET /books/_search
89 {
90   "query": {
```



```
91     "bool": {
92         "filter": [
93             {
94                 "term": {
95                     "language": "java"
96                 }
97             },
98             {
99                 "range": {
100                     "publish_time": {
101                         "gte": "2010-08-01"
102                     }
103                 }
104             }
105         ]
106     }
107 }
108 }
```

1.5 highlight高亮

highlight 关键字: 可以让符合条件的文档中的关键词高亮。

highlight相关属性:

- pre_tags 前缀标签
- post_tags 后缀标签
- tags_schema 设置为styled可以使用内置高亮样式
- require_field_match 多字段高亮需要设置为false

示例数据

```
1 #指定ik分词器
2 PUT /products
3 {
4     "settings" : {
5         "index" : {
```

```

6         "analysis.analyzer.default.type": "ik_max_word"
7     }
8 }
9 }
10
11 PUT /products/_doc/1
12 {
13     "proId" : "2",
14     "name" : "牛仔男外套",
15     "desc" : "牛仔外套男装春季衣服男春装夹克修身休闲男生潮牌工装潮流头号青年春秋棒球服男 7705浅
16     蓝常规 XL",
17     "timestamp" : 1576313264451,
18     "createTime" : "2019-12-13 12:56:56"
19 }
20
21 PUT /products/_doc/2
22 {
23     "proId" : "6",
24     "name" : "HLA海澜之家牛仔裤男",
25     "desc" : "HLA海澜之家牛仔裤男2019时尚有型舒适HKNAD3E109A 牛仔蓝(A9)175/82A(32)",
26     "timestamp" : 1576314265571,
27     "createTime" : "2019-12-18 15:56:56"
28 }

```

测试

```

1 GET /products/_search
2 {
3     "query": {
4         "term": {
5             "name": {
6                 "value": "牛仔"
7             }
8         }
9     },
10    "highlight": {
11        "fields": {
12            "*": {}
13        }
14    }
15 }

```

```
14     }  
15 }
```

自定义高亮html标签

可以在highlight中使用pre_tags和post_tags

```
1  GET /products/_search  
2  {  
3    "query": {  
4      "term": {  
5        "name": {  
6          "value": "牛仔"  
7        }  
8      }  
9    },  
10   "highlight": {  
11     "post_tags": ["</span>"],  
12     "pre_tags": ["<span style='color:red'>"],  
13     "fields": {  
14       "*": {}  
15     }  
16   }  
17 }
```

多字段高亮

```
1  
2  GET /products/_search  
3  {  
4    "query": {  
5      "term": {  
6        "name": {  
7          "value": "牛仔"  
8        }  
9      }  
10   }
```

```
9     }
10 },
11 "highlight": {
12     "pre_tags": ["<font color='red'>"],
13     "post_tags": ["<font/>"],
14     "require_field_match": "false",
15     "fields": {
16         "name": {},
17         "desc": {}
18     }
19 }
20 }
21
```

2. ES 深度分页问题及针对不同需求下的解决方案

2.1 什么是深度分页

分页问题是Elasticsearch中最常见的查询场景之一，正常情况下分页代码如实下面这样的：

```
1 # 查询第一页5条数据
2 GET /es_db/_search
3 {
4     "query": {
5         "match_all": {}
6     },
7     "from": 0,
8     "size": 5
9 }
```

输出结果如下图：

```

{
  "took" : 1,
  "timed_out" : false,
  "_shards" : {
    "total" : 1,
    "successful" : 1,
    "skipped" : 0,
    "failed" : 0
  },
  "hits" : {
    "total" : {
      "value" : 6,
      "relation" : "eq"
    },
    "max_score" : 1.0,
    "hits" : [
      { },
      { },
      { },
      { },
      { }
    ]
  ]
}

```

但是如果我们查询的数据页数特别大，当`from + size`大于10000的时候，就会出现问题，如下图报错信息所示：

ES通过参数`index.max_result_window`用来限制单次查询满足查询条件的结果窗口的大小，其默认值为10000。

2.2 深度分页会带来什么问题

ES分页查询流程大致如下：

1. 数据存储在各个分片中，协调节点将查询请求转发给各个节点，当各个节点执行搜索后，将排序后的前N条数据返回给协调节点。
2. 协调节点汇总各个分片返回的数据，再次排序，最终返回前N条数据给客户端。
3. 这个流程会导致一个深度分页的问题，也就是翻页越多，性能越差，甚至导致ES出现OOM。

在分布式系统中，对结果排序的成本随分页的深度成指数上升。

从10万名高考生中查询成绩为10001-10100位的100名考生的信息。

从上面案例中不难看出，每次有序的查询都会在每个分片中执行单独的查询，然后进行数据的二次排序，而这个二次排序的过程是发生在heap中的，也就是说当你单次查询的数量越大，那么堆内存中汇总的数据也就越多，对内存的压力也就越大。这里的单次查询的数据量取决于你查询的是第几条数据而不是查询了几条数据，比如你希望查询的是第10001-10100这一百条数据，但是ES必须将前10100

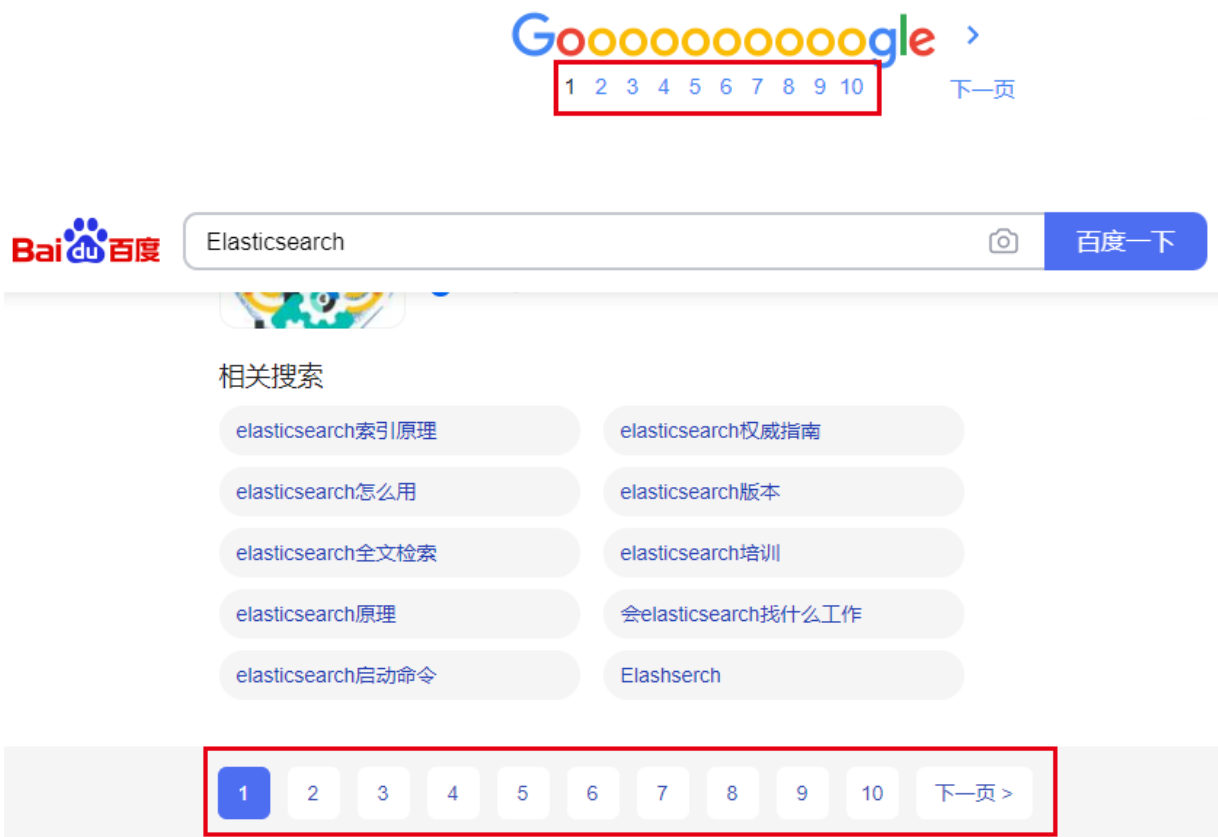
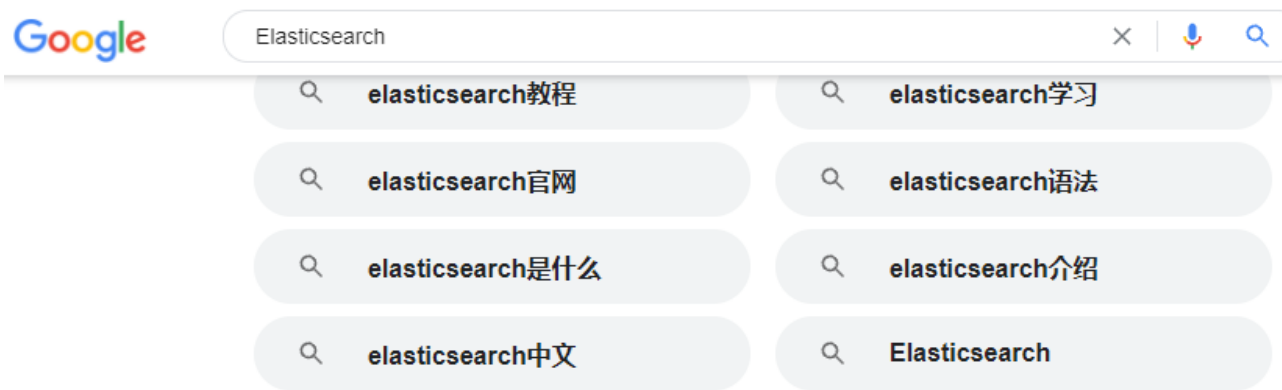
全部取出进行二次查询。因此，如果查询的数据排序越靠后，就越容易导致OOM（Out Of Memory）情况的发生，频繁的深分页查询会导致频繁的FGC。

ES为了避免用户在不了解其内部原理的情况下而做出错误的操作，设置了一个阈值，即max_result_window，其默认值为10000，其作用是为了保护堆内存不被错误操作导致溢出。

2.3 深度分页问题的常见解决方案

尝试避免使用深度分页

解决深度分页问题最好的办法就是避免使用深度分页。谷歌、百度目前作为全球和国内做大的搜索引擎不约而同的在分页条中删除了“跳页”功能，其目的就是为了避免用户使用深度分页检索。



淘宝虽然没有删除“跳页”功能，但不管我们搜索什么内容，只要商品结果足够多，返回的商品列表都是仅展示前100页的数据，其本质和ES中的max_result_window作用是一样的，都是限制你去搜索更深页数的数据。



手机端APP就更不用说了，直接是下拉加载更多，连分页条都没有，相当于你只能点击“下一页”。

滚动查询：Scroll Search

scroll滚动搜索是先搜索一批数据，然后下次再搜索下一批数据，以此类推，直到搜索出全部的数据来。

scroll搜索会在第一次搜索的时候，保存一个当时的视图快照，之后只会基于该视图快照搜索数据，如果在搜索期间数据发生了变更，用户是看不到变更的数据的。因此，滚动查询不适合实时性要求高的搜索场景。

官方已不推荐使用滚动查询进行深度分页查询，因为无法保存索引状态。

适合场景

单个滚动搜索请求中检索大量结果，即非“C端业务”场景

使用

1) 第一次进行scroll查询：

```
1 #查询命令中新增scroll=1m,说明采用游标查询，保持游标查询窗口1分钟，也就是本次快照的结果缓存起来的有效时间是1分钟。
2 #这里由于测试数据量不够，所以size值设置为2。
3 #实际使用中为了减少游标查询的次数，可以将值适当增大，比如设置为1000。
4 GET /es_db/_search?scroll=1m
5 {
6   "query": { "match_all": {}},
7   "size": 2
8 }
```

查询结果：除了返回前2条记录，还返回了一个游标ID值_scroll_id。

2) 从第二次查询开始，每次查询都要指定_scroll_id参数：

```
1 # scroll_id 的值就是上一个请求中返回的 _scroll_id 的值
2 GET /_search/scroll
3 {
4     "scroll": "1m",
5     "scroll_id" :
6     "FglUy2x1ZGVfY29udGV4dF91dWlkDXF1ZXJ5QW5kRmV0Y2gBFmNwcVdjblRxUzVhZXlicG9HeU02bWcAAAAAAABmzRY2YlV3Z0o5VVNTdwJobkE5Z3MtXzJB"
```

查询结果:

多次根据scroll_id游标查询，直到没有数据返回则结束查询。采用游标查询索引全量数据，更安全高效，限制了单次对内存的消耗。

删除游标scroll

scroll超过超时后，搜索上下文会自动删除。然而，保持scroll打开是有代价的，因此一旦不再使用，就应明确清除scroll上下文

```
1 DELETE /_search/scroll
2 {
3     "scroll_id" :
4     "FGluY2x1ZGVfY29udGV4dF91dWlkDXF1ZXJ5QW5kRmV0Y2gBFmNwcVdjblRxUzVhZX1icG9HeU02bWcAAAAAABmzRY2YlV3Z0o5VVNTdWJ0bkE5Z3MtXzJB"
```

注意事项

- scroll滚动查询不适合实时性要求高的查询场景，比较适合数据迁移的场景。
- scroll查询完毕后，要手动清理掉 scroll_id。虽然ES有自动清理机制，但是 scroll_id 的存在会耗费大量的资源来保存一份当前查询结果集映像，并且会占用文件描述符。

官方建议：ES7之后，不再建议使用scroll API进行深度分页。如果要分页检索超过 Top 10,000+ 结果时，推荐使用：PIT + search after。

search after

参考文档: <https://www.elastic.co/guide/en/elasticsearch/reference/7.17/paginate-search-results.html#search-after>

scroll API适用于高效的深度滚动，但滚动上下文成本高昂，不建议将其用于实时用户请求。而search_after参数通过提供一个活动光标来规避这个问题。这样可以使用上一页的结果来帮助检索下一页。

search_after 分页查询可以简单概括为如下几个步骤：

1) 获取索引的pit

使用 search_after 需要具有相同查询和排序值的多个搜索请求。如果在这些请求之间发生刷新，结果的顺序可能会发生变化，从而导致跨页面的结果不一致。为防止出现这种情况，可以创建一个时间点 (PIT) 以保留搜索中的当前索引状态。Point In Time (PIT) 是 Elasticsearch 7.10 版本之后才有的新特性。

```
1 # 创建一个时间点(PIT)来保存搜索期间的当前索引状态
2 POST /es_db/_pit?keep_alive=1m
3 #返回结果，会返回一个PID的值
4 {
5   "id" :
6     "39K1AwEFZXNfZGIWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAWdkhjbE9YNVRTMUNDcWNQQVR2ZXYzdAAAAAAAAA
7     A9jhZvaGpLSDlzMVMxbW5idG5DZ0xEUHFRAAEWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAA"
```

2) 根据pit首次查询

根据pit查询的时候，不用指定索引的名词

```
1 GET /_search
2 {
3   "query": {
4     "match_all": {}
5   },
6   "pit": {
7     "id":
8       "39K1AwEFZXNfZGIWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAWdkhjbE9YNVRTMUNDcWNQQVR2ZXYzdAAAAAAAAA
9       A9jhZvaGpLSDlzMVMxbW5idG5DZ0xEUHFRAAEWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAA",
10     "keep_alive": "1m"
11   },
12   "size": 2,
13   "sort": [
14     {"_id": "asc"}
15   ]
16 }
```

返回结果:

```
1  {
2    "pit_id" :
3      "39K1AwEFZXNfZGIWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAwdkhjbE9YNVRTMUNDcWNQQVR2ZXYzdwAAAAAAA
4      A7hRZvaGpLSDlzVVMxbW5idG5DZ0xEUHFRAAEWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAA",
5    "took" : 16,
6    "timed_out" : false,
7    "_shards" : {
8      "total" : 1,
9      "successful" : 1,
10     "skipped" : 0,
11     "failed" : 0
12   },
13   "hits" : {
14     "total" : {
15       "value" : 5,
16       "relation" : "eq"
17     },
18     "max_score" : null,
19     "hits" : [
20       {
21         "_index" : "es_db",
22         "_type" : "_doc",
23         "_id" : "2",
24         "_score" : null,
25         "_source" : {
26           "name" : "李四",
27           "sex" : 1,
28           "age" : 28,
29           "address" : "广州荔湾大厦",
30           "remark" : "java assistant"
31         },
32         "sort" : [
33           "2",
34           0
35         ]
36       },
37     ]
38   },
39   {
```

```

36         "_index" : "es_db",
37         "_type" : "_doc",
38         "_id" : "3",
39         "_score" : null,
40         "_source" : {
41             "name" : "王五",
42             "sex" : 0,
43             "age" : 26,
44             "address" : "广州白云山公园",
45             "remark" : "php developer"
46         },
47         "sort" : [
48             "3",
49             1
50         ]
51     }
52 ]
53 }
54 }

```

3) 根据search_after和pit进行翻页查询

要获得下一页结果，请使用最后一次命中的排序值（包括 tiebreaker）作为 search_after 参数重新运行先前的搜索。 如果使用 PIT，请在 pit.id 参数中使用最新的 PIT ID。 搜索的查询和排序参数必须保持不变。

```

1  #search_after指定为上一次查询返回的sort值。
2  GET /_search
3  {
4      "query": {
5          "match_all": {}
6      },
7      "pit": {
8          "id":
9              "39K1AwEFZXNfZGIWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAWdkhjbE9YNVRTMUNDcWNQQVR2ZXYzdWAAAAAAAAA
10             A9jhZvaGpLSDlzMVMxbW5idG5DZ0xEUHFRAAEWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAA",
11          "keep_alive": "1m"
12      },
13      "size": 2,
14      "sort": [

```

```
13         {"_id": "asc"}
14     ],
15     "search_after": [
16         5
17     ]
18 }
```

返回结果:

```
1  {
2    "pit_id" :
3    "39K1AwEFZXNfZGIWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAWdkhjbE9YNVRTMUNDcWNQQVR2ZXYzdwAAAAAAAAA
4    A8wxZvaGpLSDlzMVMxbW5idG5DZ0xEUHFRAAEWZTN2N2Nrdk5RRjY3QjBma1h5aFRodwAA",
5    "took" : 1,
6    "timed_out" : false,
7    "_shards" : {
8      "total" : 1,
9      "successful" : 1,
10     "skipped" : 0,
11     "failed" : 0
12   },
13   "hits" : {
14     "total" : {
15       "value" : 5,
16       "relation" : "eq"
17     },
18     "max_score" : null,
19     "hits" : [
20       {
21         "_index" : "es_db",
22         "_type" : "_doc",
23         "_id" : "4",
24         "_score" : null,
25         "_source" : {
26           "name" : "赵六",
27           "sex" : 0,
28           "age" : 22,
29           "address" : "长沙橘子洲",
30           "remark" : "python assistant"
```

```
29     },
30     "sort" : [
31         "4"
32     ]
33 },
34 {
35     "_index" : "es_db",
36     "_type" : "_doc",
37     "_id" : "5",
38     "_score" : null,
39     "_source" : {
40         "name" : "张龙",
41         "sex" : 0,
42         "age" : 19,
43         "address" : "长沙麓谷企业广场",
44         "remark" : "java architect assistant"
45     },
46     "sort" : [
47         "5"
48     ]
49 }
50 ]
51 }
52 }
```

2.4 总结

分页方式	性能	优点	缺点	适用场景
from + size	低	灵活性好，实现简单，支持随机翻页	受制于max_result_window设置，不能无限制翻页；存在深度翻译问题，越往后翻译越慢。	数据量比较小，能容忍深度分页问题
scroll	中	解决了深度分页问题	scroll查询的相应数据是非实时的，如果遍历过程中插	海量数据的导出，需要查询海量结果集的数据

			入新的数据，是查询不到的； 保留上下文需要足够的堆内存空间。	
search_after	高	性能最好，不存在深度分页问题，能够反映数据的实时变更	实现复杂，需要有一个全局唯一的字段连续分页的实现会比较复杂，因为每一次查询都需要上次查询的结果，它不适用于大幅度跳页查询	海量数据的分页