

## 1. 聚合的概述

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。

## 2. 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}
```

## 指标聚合 (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 }
```

## 桶聚合 (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

桶聚合可以用于各种场景，例如：

- 对数据进行分组统计，比如按照地区、年龄段、性别等字段进行分组统计。
- 对时间序列数据进行时间段分析，比如按照每小时、每天、每月、每季度、每年等时间段进行分析。
- 对各种标签信息分类，并统计其数量。

## 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  }
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 }
46 }
```

## 管道聚合 (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 }
```

```

40     }
41   }},
42   "all_jobs": {
43     "terms": {
44       "field": "job.keyword"
45     }
46   }
47 }
48 }
49 }

```

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 }

```

71 **#global**

72 **#** 使用**global**聚合来计算所有匹配查询的文档（即所有年龄大于或等于**40**岁的员工）的平均薪资。

73 **#global**聚合的特点是它会考虑查询范围内的所有文档，而不仅仅是某个特定分组或桶中的文档。

74 **POST** employees/\_search

```

75 {
76   "size": 0,
77   "query": {
78     "range": {
79       "age": {

```

```
80         "gte": 40
81     }
82 }
83 },
84 "aggs": {
85     "jobs": {
86         "terms": {
87             "field": "job.keyword"
88         }
89     },
90
91     "all": {
92         "global": {},
93         "aggs": {
94             "salary_avg": {
95                 "avg": {
96                     "field": "salary"
97                 }
98             }
99         }
100     }
101 }
102 }
103
104
105
```

## 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 }
```

```
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             "aggs": {
69                 "stats_salary": {
70                     "stats": {
71                         "field": "salary"
72                     }
73                 }
74             }
75         }
76     }
77 }
```

### 3. 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. Elasticsearch 聚合性能优化

### 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。你可以执行一个带有过滤查询的搜索请求，Elasticsearch 将自动尝试使用节点查询缓存来优化性能。例如，如果你想缓存一个基于特定字段值的过滤查询，你可以发送如下的 HTTP 请求：



```
1 GET /your_index/_search
2 {
3   "query": {
4     "bool": {
5       "filter": {
6         "term": {
7           "your_field": "your_value"
8         }
9       }
10    }
11  }
12 }
```

## 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"}}}}
```