# Elasticsearch Database Report

Junhui LI, Ge QIU, Zhiheng WU, Yue JI

Group : A4 IBO4

**Dataset :**

Stock exchange

**Abstract :**

Our group uses the data set named stock exchange, each item contains attributes such as id, compamy, price, date, etc. We design related queries based on these attributes.

## 1. Introduction

Elasticsearch is an open source search engine highly scalable. It allows you to keep and analyse a great volume of information practically in real time.

Elasticsearch works with JSON documents files. Using an internal structure, it can parse your data in almost real time to search for the information you need.

## 2. Data prepossessing

Elasticsearch has features like *index* and *_type*. So we need to process the data before inserting the data. We add types and indexes in bulk in the dataset. Here is the program python for adding.

```python
import json
import os

DATASET_FILE = 'stocks.json'
BASE_FILE = 'movielens_usersRating.json'
TYPENAME = 'stock'
INDEXNAME = 'stocks'

if not os.path.isfile(DATASET_FILE):
    print("[#] starting json modification...")
    with open(BASE_FILE, "r") as f:
        with open(DATASET_FILE, "a+") as f2:
            for line in f:
                data = json.loads(line.strip('\n'))
                print data

                index = {"index":{"_index":INDEXNAME, "_type":TYPENAME,
"_id":data["_id"]}}
                data.pop("_id")
                fields={}
                fields["fields"] = data
```

```
                    # data["field"] = data["_id"]["$oid"]
                    f2.write(json.dumps(index)+'\n')
                    f2.write(json.dumps(fields)+'\n')
        print("[#] modification done")
```

### 3. Import dataset with cURL

```
curl -XPUT localhost:9200/_bulk -H"Content-Type: application/json" --data-
binary @stock.json
```
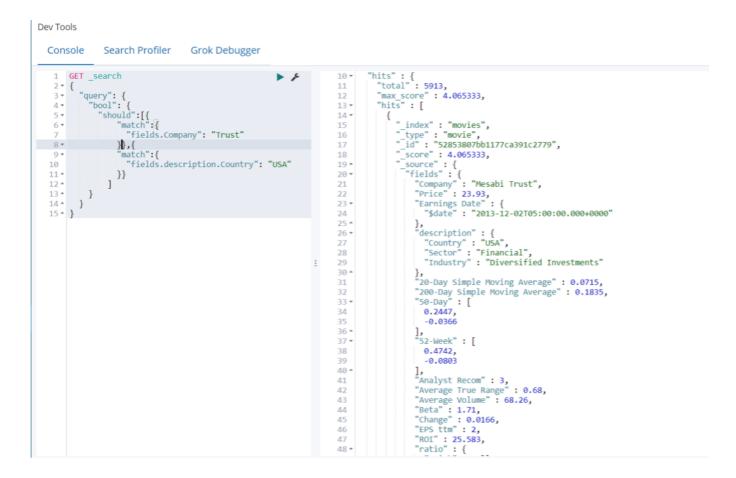
### 4. Queries

- Simple queries:

    *1. find the data of company name which contains "Trust"(simple match)*

```
GET _search
{
  "query": {
    "match" : {
      "fields.Company": "Trust"
    }
  }
}
```

2. *match the company contains "Trust" and Country contains "USA"(simple should)*

```
GET _search
{
  "query": {
    "bool": {
      "should":[{
          "match":{
            "fields.Company": "Trust"
          }},{
          "match":{
            "fields.description.Country": "USA"
          }}
        ]
      }
    }
}
```

```
Dev Tools

  Console    Search Profiler    Grok Debugger

 1   GET _search                    ▶ ⚙        10▾   "hits" : {
 2▾  {                                          11      "total" : 5913,
 3▾    "query": {                               12      "max_score" : 4.065333,
 4▾      "bool": {                              13▾     "hits" : [
 5▾        "should":[{                          14▾       {
 6▾          "match":{                          15          "_index" : "movies",
 7            "fields.Company": "Trust"         16          "_type" : "movie",
 8▾          }},{                               17          "_id" : "52853807bb1177ca391c2779",
 9▾          "match":{                          18          "_score" : 4.065333,
10            "fields.description.Country": "USA"  19▾       "_source" : {
11▴          }}                                 20▾          "fields" : {
12▴        ]                                    21            "Company" : "Mesabi Trust",
13▴      }                                      22            "Price" : 23.93,
14▴    }                                        23▾           "Earnings Date" : {
15▴  }                                          24              "$date" : "2013-12-02T05:00:00.000+0000"
                                                25▴           },
                                                26▾           "description" : {
                                                27              "Country" : "USA",
                                                28              "Sector" : "Financial",
                                                29              "Industry" : "Diversified Investments"
                                                30▴           },
                                                31            "20-Day Simple Moving Average" : 0.0715,
                                                32            "200-Day Simple Moving Average" : 0.1835,
                                                33▾           "50-Day" : [
                                                34              0.2447,
                                                35              -0.0366
                                                36▴           ],
                                                37▾           "52-Week" : [
                                                38              0.4742,
                                                39              -0.0803
                                                40▴           ],
                                                41            "Analyst Recom" : 3,
                                                42            "Average True Range" : 0.68,
                                                43            "Average Volume" : 68.26,
                                                44            "Beta" : 1.71,
                                                45            "Change" : 0.0166,
                                                46            "EPS ttm" : 2,
                                                47            "ROI" : 25.583,
                                                48▾           "ratio" : {
```

*3. find that the Sector must be Services and Company should contains "America"*

```
GET _search
{
  "query": {
    "bool": {
      "should":{
        "match":{
          "fields.Company": "America"
        }
      },
      "must":{
        "match":{
          "fields.description.Sector": "Services"
        }
      }
    }
  }
}
```

Console     Search Profiler     Grok Debugger

```
1  GET _search                          ▶ ⚒
2 ▾ {
3 ▾   "query": {
4 ▾     "bool": {
5 ▾       "should":{
6 ▾         "match":{
7             "fields.Company": "America"
8 ▴         }
9 ▴       },
10 ▾       "must":{
11 ▾         "match":{
12             "fields.description.Sector": "Services"
13 ▴         }
14 ▴       }
15 ▴     }
16 ▴   }
17 ▴ }
18   |
```

```
1 ▾ {
2      "took" : 0,
3      "timed_out" : false,
4 ▾   "_shards" : {
5        "total" : 6,
6        "successful" : 6,
7        "skipped" : 0,
8        "failed" : 0
9 ▴   },
10 ▾   "hits" : {
11       "total" : 864,
12       "max_score" : 8.124348,
13 ▾     "hits" : [
14 ▾       {
15           "_index" : "movies",
16           "_type" : "movie",
17           "_id" : "52853802bb1177ca391c1b4b",
18           "_score" : 8.124348,
19 ▾         "_source" : {
20 ▾           "fields" : {
21               "Company" : "Bowl America Inc.",
22               "Price" : 15.07,
23               "Earnings Date" : null,
24 ▾             "description" : {
25                 "Country" : "USA",
26                 "Sector" : "Services",
27                 "Industry" : "Sporting Activities"
28 ▴             },
29               "20-Day Simple Moving Average" : 0.0433,
30               "200-Day Simple Moving Average" : 0.1548,
31 ▾             "50-Day" : [
32                 0.1339,
33                 -0.0118
34 ▴             ],
35 ▾             "52-Week" : [
36                 0.3486,
37                 -0.0118
38 ▴             ],
39               "Analyst Recom" : 3,
```

*4. find the Industry named "Management Services"(simple match_phrase)*

```
GET _search
{
  "query": {
    "match_phrase": {
      "fields.description.Industry": "Management Services"
    }
  }
}
```

Dev Tools

Console        Search Profiler        Grok Debugger

```
1  GET _search                                    ▶ 🔧
2  {
3    "query": {
4      "match_phrase": {
5        "fields.description.Industry": "Management Services"
6      }
7    }
8  }
9
```

```
10  "hits" : {
11    "total" : 21,
12    "max_score" : 6.9696426,
13    "hits" : [
14      {
15        "_index" : "movies",
16        "_type" : "movie",
17        "_id" : "52853804bb1177ca391c2035",
18        "_score" : 6.9696426,
19        "_source" : {
20          "fields" : {
21            "Company" : "Exponent Inc.",
22            "Price" : 74.51,
23            "Earnings Date" : {
24              "$date" : "2013-10-16T20:30:00.000+0000"
25            },
26            "description" : {
27              "Country" : "USA",
28              "Sector" : "Services",
29              "Industry" : "Management Services"
30            },
31            "20-Day Simple Moving Average" : -0.0239,
32            "200-Day Simple Moving Average" : 0.224,
33            "50-Day" : [
34              0.1481,
35              -0.0744
36            ],
37            "52-Week" : [
38              0.5919,
39              -0.0744
40            ],
41            "Analyst Recom" : 2.5,
42            "Average True Range" : 1.65,
43            "Average Volume" : 33.05,
44            "Beta" : 0.78,
45            "Change" : 0.0016,
46            "EPS ttm" : 2.72,
47            "ROI" : 0.153,
48            "ratio" : {
```

*5. find the ROI greater than 2(simple range)*

```
GET _search
{
  "query": {
    "bool":{
      "must":{
        "range":{
          "fields.ROI":{
            "gte": 2
          }
        }
      }
    }
  }
}
```

Console    Search Profiler    Grok Debugger

```
1  GET _search                                    ▶ 🔧
2  {
3    "query": {
4      "bool":{
5        "must":{
6          "range":{
7            "fields.ROI":{
8              "gte": 2
9            }
10         }
11       }
12     }
13   }
14  }
15
```

```
10  "hits" : {
11    "total" : 21,
12    "max_score" : 1.0,
13    "hits" : [
14      {
15        "_index" : "movies",
16        "_type" : "movie",
17        "_id" : "52853806bb1177ca391c24e5",
18        "_score" : 1.0,
19        "_source" : {
20          "fields" : {
21            "Company" : "James Hardie Industries SE",
22            "Price" : 54.25,
23            "Earnings Date" : {
24              "$date" : "2013-11-14T05:00:00.000+0000"
25            },
26            "description" : {
27              "Country" : "Ireland",
28              "Sector" : "Industrial Goods",
29              "Industry" : "Cement"
30            },
31            "20-Day Simple Moving Average" : 0.0646,
32            "200-Day Simple Moving Average" : 0.1466,
33            "50-Day" : [
34              0.2363,
35              0.0294
36            ],
37            "52-Week" : [
38              0.3573,
39              0.0294
40            ],
41            "Analyst Recom" : 1,
42            "Average True Range" : 0.8,
43            "Average Volume" : 3.59,
44            "Beta" : 1.98,
45            "Change" : 0.1142,
46            "EPS ttm" : 1.34,
47            "ROI" : 2.269,
48            "ratio" : {
```

**6. find the Earning Date that between 2014-01-01 and 2015-01-01(simple date)**

```
GET _search
{
  "query": {
    "bool":{
      "must":{
        "range":{
          "fields.Earnings Date.$date":{
            "from": "2014-01-01",
            "to": "2015-01-01"
          }
        }
      }
    }
  }
}
```

- Complex queries:

  *1. show the year performance aggregation in desc order(simple terms)*

```
GET _search
{
  "aggs": {
    "years": {
      "terms": {
        "field": "fields.performance.Year",
        "size": 10,
        "order": {
          "_count": "desc"
        }
      }
    }
  }
}
```

Console    Search Profiler    Grok Debugger

```
1  GET _search                                    ▶ ⚙
2  {
3    "aggs": {
4      "years": {
5        "terms": {
6          "field": "fields.performance.Year",
7          "size": 10,
8          "order": {
9            "_count": "desc"
10         }
11       }
12     }
13   }
14 }
15
```

```
353  "aggregations" : {
354    "years" : {
355      "doc_count_error_upper_bound" : 10,
356      "sum_other_doc_count" : 6365,
357      "buckets" : [
358        {
359          "key" : 0.0,
360          "doc_count" : 8
361        },
362        {
363          "key" : 0.06920000165700912,
364          "doc_count" : 5
365        },
366        {
367          "key" : -0.19910000264644623,
368          "doc_count" : 3
369        },
370        {
371          "key" : -0.02969999983906746,
372          "doc_count" : 3
373        },
374        {
375          "key" : -0.018200000748038292,
376          "doc_count" : 3
377        },
378        {
379          "key" : -0.003000000026077032,
380          "doc_count" : 3
381        },
382        {
383          "key" : 0.0015999999595806003,
384          "doc_count" : 3
385        },
386        {
387          "key" : 0.04699999839067459,
388          "doc_count" : 3
389        },
```

*2. Group by performance.year and find their average ROI(composite)*

```
GET _search
{
  "aggs": {
    "group_by_year": {
      "terms": {
        "field": "fields.performance.Year"
      },
      "aggs": {
        "avg_by_ROI": {
          "avg": {
            "field": "fields.ROI"
          }
        }
      }
    }
  }
}
```

Console    Search Profiler    Grok Debugger



### 3. number of distinct industry(cardinality)



So first put the industry as fielddata.

```
GET _search
{
  "aggs": {
    "distinct_count": {
      "cardinality": {
        "field": "fields.description.Industry"
      }
    }
  }
}
PUT /movies/movie/_mappings
```

```
{ "properties": { "fields.description.Industry": { "type": "text",
"fielddata": true } } }
```

Console    Search Profiler    Grok Debugger

```
 1  GET _search                              ▶ 🔧     323 ▾        "50-Day" : [
 2 ▾ {                                               324            0.123,
 3 ▾   "aggs": {                                     325            -0.0574
 4 ▾     "distinct_count": {                         326 ▴        ],
 5 ▾       "cardinality": {                          327 ▾        "52-Week" : [
 6           "field": "fields.description.Industry"  328            0.4687,
 7 ▴       }                                         329            -0.1859
 8 ▴     }                                           330 ▴        ],
 9 ▴   }                                             331          "Analyst Recom" : 3,
10 ▴ }                                               332          "Average True Range" : 0.59,
11  PUT /movies/movie/_mappings                      333          "Average Volume" : 348.08,
12  { "properties": { "fields.description.Industry": { "type":  334          "Beta" : 1.63,
    "text", "fielddata": true } } }                  335          "Change" : -9.0E-4,
13                                                    336          "EPS ttm" : 1.36,
14                                                    337          "ROI" : 0.033,
                                                      338 ▾        "ratio" : {
                                                      339            "quick" : null,
                                                      340            "current" : null
                                               ⋮      341 ▴        },
                                                      342 ▾        "performance" : {
                                                      343            "Year" : 0.3884,
                                                      344            "Half Year" : -0.1159,
                                                      345            "Month" : 0.0847,
                                                      346            "Week" : -0.0302
                                                      347 ▴        }
                                                      348 ▴      }
                                                      349 ▴    }
                                                      350 ▴  }
                                                      351 ▴  ]
                                                      352 ▴  },
                                                      353 ▾  "aggregations" : {
                                                      354 ▾    "distinct_count" : {
                                                      355        "value" : 268
                                                      356 ▴    }
                                                      357 ▴  }
                                                      358 ▴ }
                                                      359
```

- Hard query:

  *filter the average volume that greater than 5000, then group by the performance of the year, finally calculate their average ROI*

```
GET _search
{
    "query" : {
        "bool" : {
          "must" : {
             "range" : {
                "fields.Average Volume":{
                   "gte" : 5000
                }
             }
          }
        }
    },
    "aggs": {
      "group_by_year": {
        "terms": {
          "field": "fields.description.Country",
          "size": 5,
```

```
          "order": {
            "_count": "desc"
          }
        },
        "aggs": {
          "average_of_ROI": {
            "avg": {
              "field": "fields.ROI"
            }
          }
        }
      }
    }
  }
}
PUT /movies/movie/_mappings
{ "properties": { "fields.description.Country": { "type": "text",
"fielddata": true } } }
```



## 5. Disscussion and Conclusion

- Problems： Group by range

Dev Tools                                                                    History    Settings    Hel

Console        Search Profiler        Grok Debugger

```
 1   GET _search                          ▶ 🔧        1 ▾ {
 2 ▾ {                                               2 ▾     "error": {
 3 ▾     "aggs": {                                   3 ▾       "root_cause": [
 4 ▾       "range_age": {                            4 ▾         {
 5 ▾         "terms": {                              5             "type": "x_content_parse_exception",
 6             "field": "fields.performance.Year",   6             "reason": "[6:9] [terms] unknown field [ranges],
 7 ▾           "ranges": [                                          parser not found"
 8 ▾             {                                   7 ▴         }
 9               "from": -1,                         8 ▴       ],
10               "to": 0                             9         "type": "x_content_parse_exception",
11 ▴           }                                    10         "reason": "[6:9] [terms] unknown field [ranges], parser
12 ▴         ]                                                    not found"
13 ▴       }                                        11 ▴     },
14 ▴     }                                          12       "status": 400
15 ▴   }                                            13 ▴ }
16 ▴ }
17
```

Cannot find proper position for the ranges within a single aggs. But we can do it by mapping or composite aggregations.

- Conclusion

This project, given the many difficulties encountered, really allowed us to understand the basics of Elasticsearch. We realized that an organization, clarity and good structure in the Elasticsearch is paramount so as not to fall into the mistakes and pitfalls of such language.

Finally, this project allowed us to concretely practice the Elasticsearch and the data structure, which may later in our journey help us better understand the logic of NoSQL.