

به نام خدا



دانشگاه تهران

پردیس دانشکده‌های فنی

دانشکده برق و کامپیوتر

آزمایشگاه پایگاه داده

دستور کار شماره ۷

نام و نام خانوادگی

معین شیردل ۸۱۰۱۹۷۵۳۵

آبان ماه ۱۴۰۰

بخش اول) آموزش های سایت bigdata.ir:

```
2 * {
3   "title": "The Godfather",
4   "director": "Francis Ford Coppola",
5   "year": 1972
6 * }
7
8 PUT movies/_doc/2
9 * {
10  "title": "Lawrence of Arabia",
11  "director": "David Lean",
12  "year": 1962,
13  "genres": ["Adventure", "Biography", "Drama"]
14 * }
15 PUT movies/_doc/3
16 * {
17  "title": "To Kill a Mockingbird",
18  "director": "Robert Mulligan",
19  "year": 1962,
20  "genres": ["Crime", "Drama", "Mystery"]
21 * }
22 PUT movies/_doc/4
23 * {
24  "title": "Apocalypse Now",
25  "director": "Francis Ford Coppola",
26  "year": 1979,
27  "genres": ["Drama", "War"]
28 * }
29 PUT movies/_doc/5
30 * {
31  "title": "Kill Bill: Vol. 1",
32  "director": "Quentin Tarantino",
33  "year": 2003,
34  "genres": ["Action", "Crime", "Thriller"]
35 * }
36 PUT movies/_doc/6
37 * {
38  "title": "The Assassination of Jesse James by the Coward Robert Ford",
39  "director": "Andrew Dominik",
40  "year": 2007,
41  "genres": ["Biography", "Crime", "Drama"]
42 * }
```

```
1 #! Elasticsearch built-in security features are not enabled. Without authentication, your
2 cluster could be accessible to anyone. See https://www.elastic.co/guide/en/elasticsearch
3 /reference/7.16/security-minimal-setup.html to enable security.
4
5 {
6   "_index" : "movies",
7   "_type" : "_doc",
8   "_id" : "6",
9   "_version" : 2,
10  "result" : "updated",
11  "shards" : {
12    "total" : 2,
13    "successful" : 1,
14    "failed" : 0
15  },
16  "_seq_no" : 7,
17  "_primary_term" : 1
18 }
```

به کمک دستورات بالا تعدادی رکورد در دیتابیس محلی elasticsearch ذخیره می کنیم تا داده برای شروع داشته باشیم. خروجی، حاصل اجرای آخرین دستور است.

```
3
4
5 GET /movies/_doc/1
6
7
```

```
2 {
3   "_index" : "movies",
4   "_type" : "_doc",
5   "_id" : "1",
6   "_version" : 2,
7   "_seq_no" : 1,
8   "_primary_term" : 1,
9   "found" : true,
10  "_source" : {
11    "title" : "The Godfather",
12    "director" : "Francis Ford Coppola",
13    "year" : 1972
14  }
15 }
16
```

به کمک دستور GET به صورت بالا می توانیم اطلاعات رکورد (یا در این مثال فیلم) مورد نظر که ID آن را می دانیم را بدست آوریم. در حالت بالا اطلاعات فیلم اول را میبینیم.

```
POST _all/_search
{
  "query": {
    "query_string": {
      "query": "kill"
    }
  }
}

3 {
4   "took" : 23,
5   "timed_out" : false,
6   "_shards" : {
7     "total" : 5,
8     "successful" : 5,
9     "skipped" : 0,
10    "failed" : 0
11  },
12  "hits" : {
13    "total" : {
14      "value" : 2,
15      "relation" : "eq"
16    },
17    "max_score" : 1.2667978,
18    "hits" : [
19      {
20        "_index" : "movies",
21        "_type" : "movie",
22        "_id" : "3",
23        "_score" : 1.2667978,
24        "_source" : {
25          "title" : "To Kill a Mockingbird",
26          "director" : "Robert Mulligan",
27          "year" : 1962,
28          "genres" : [
29            "Crime",
30            "Drama",
31            "Mystery"
32          ]
33        }
34      },
35      {
36        "_index" : "movies",
37        "_type" : "movie",
38        "_id" : "5",
39        "_score" : 1.2667978,
40        "_source" : {
41          "title" : "Kill Bill: Vol. 1",
42          "director" : "Quentin Tarantino",
43          "year" : 2003
44        }
45      }
46    ]
47  }
48 }
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
```

با جستجو در تمامی رکوردها، رکوردهایی که عبارت Kill در آن ها به کار رفته است نمایش داده شدند که دو رکورد مشاهده شده به واسطه نامشان که Kill Bill و To Kill a Mockingbird است در نتایج آمده اند.

```
GET /movies/_search
{
  "query": {
    "query_string": {
      "query": "title:ford OR kill"
    }
  }
}

29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69

{
  "index" : "movies",
  "type" : "movie",
  "id" : "5",
  "score" : 1.2667978,
  "source" : {
    "title" : "Kill Bill: Vol. 1",
    "director" : "Quentin Tarantino",
    "year" : 2003,
    "genres" : [
      "Action",
      "Crime",
      "Thriller"
    ]
  }
},
{
  "index" : "movies",
  "type" : "movie",
  "id" : "6",
  "score" : 0.82546186,
  "source" : {
    "title" : "The Assassination of Jesse James by the Coward Robert Ford",
    "director" : "Andrew Dominik",
    "year" : 2007,
    "genres" : [
      "Biography",
      "Crime",
      "Drama"
    ]
  }
}
}
```

```
GET /movies/_search
{
  "query": {
    "query_string": {
      "query": "(title:kill AND (director:Tarantino))"
    }
  }
}

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" : 1,
14      "relation" : "eq"
15    },
16    "max_score" : 3.0309887,
17    "hits" : [
18      {
19        "_index" : "movies",
20        "_type" : "movie",
21        "_id" : "5",
22        "_score" : 3.0309887,
23        "_source" : {
24          "title" : "Kill Bill: Vol. 1",
25          "director" : "Quentin Tarantino",
26          "year" : 2003,
27          "genres" : [
28            "Action",
29            "Crime",
30            "Thriller"
31          ]
32        }
33      }
34    ]
35  }
36 }
37
```

به شکل های مختلف می توان شروط متفاوت روی مقادیر فیلد ها گذاشت. مثلا در تصویر اول نتایجی که در title خود کلمه Ford را دارند را نیز خواسته ایم که یک رکورد به رکورد های قبلی اضافه کرده است و در تصویر دوم، با شروط منطقی مشخص کرده ایم که رکورد هایی با title دارای عبارت Ford و director دارای کلمه Tarantino را می خواهیم.

```
POST /movies/_search
{
  "query": {
    "match": {
      "genres": {
        "query": "drama"
      }
    }
  }
}
```

```
2 * {
3   "took" : 0,
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" : 4,
14      "relation" : "eq"
15    },
16    "max_score" : 0.27414778,
17    "hits" : [
18      {
19        "_index" : "movies",
20        "_type" : "movie",
21        "_id" : "4",
22        "_score" : 0.27414778,
23        "_source" : {
24          "title" : "Apocalypse Now",
25          "director" : "Francis Ford Coppola",
26          "year" : 1979,
27          "genres" : [
28            "Drama",
29            "War"
30          ]
31        }
32      },
33      {
34        "_index" : "movies",
35        "_type" : "movie",
36        "_id" : "2",
37        "_score" : 0.23549506,
38        "_source" : {
39          "title" : "Lawrence of Arabia",
40          "director" : "David Lean",
41          "year" : 1962,
42          "genres" : [
```

به کمک دستور بالا فیلم هایی با ژانر درام را پیدا میکنیم. این ساختار، ساختار عمومی تر برای اعمال فیلترها روی فیلدهاست.

```
POST /movies/_search
{
  "query": {
    "match": {
      "title": {
        "query": "Kill Bill",
        "operator": "or"
      }
    }
  }
}
```

```
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" : 2,
14      "relation" : "eq"
15    },
16    "max_score" : 3.0899403,
17    "hits" : [
18      {
19        "_index" : "movies",
20        "_type" : "movie",
21        "_id" : "5",
22        "_score" : 3.0899403,
23        "_source" : {
24          "title" : "Kill Bill: Vol. 1",
25          "director" : "Quentin Tarantino",
26          "year" : 2003,
27          "genres" : [
28            "Action",
29            "Crime",
30            "Thriller"
31          ]
32      },
33      {
34        "_index" : "movies",
35        "_type" : "movie",
36        "_id" : "3",
37        "_score" : 1.2667978,
38        "_source" : {
39          "title" : "To Kill a Mockingbird",
40          "director" : "Robert Mulligan",
41          "year" : 1962,
42          "genres" : [
```

هنگامی که کلمه query مانند Kill Bill باشد، الستیک آن را به مانند دو کلمه برای سرچ می بیند و به دنبال میچ کردن هر دو کلمه می گردد. Operator نشان دهنده نحوه قرار گرفتن نتایج هر دو میچ کردن در جواب نهاییست. می بینیم که حاصل این کوئری دو تا فیلم است که فیلم با عنوان Kill Bill Vol. 1 یا هر دو کلمه میچ شده و score بالاتری هم گرفته است.

```
POST /movies/_search
{
  "query": {
    "multi_match": {
      "query": "ford",
      "fields": [
        "titles",
        "director"
      ]
    }
  }
}
```

```
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" : 3,
14      "relation" : "eq"
15    },
16    "max_score" : 2.4763856,
17    "hits" : [
18      {
19        "_index" : "movies",
20        "_type" : "movie",
21        "_id" : "6",
22        "_score" : 2.4763856,
23        "_source" : {
24          "title" : "The Assassination of Jesse James by the Coward Robert Ford",
25          "director" : "Andrew Dominik",
26          "year" : 2007,
27          "genres" : [
28            "Biography",
29            "Crime",
30            "Drama"
31          ]
32        }
33      },
34      {
35        "_index" : "movies",
36        "_type" : "movie",
37        "_id" : "1",
38        "_score" : 1.0313075,
39        "_source" : {
40          "title" : "The Godfather",
41          "director" : "Francis Ford Coppola",
42          "year" : 1972
43        }
44      }
45    ]
46  }
47 }
```

برای میچ کردن روی دو فیلد مختلف هم از multi_match استفاده می کنیم. در اینجا کلمه Ford را روی هر دو فیلد عنوان و کارگردان میچ کرده و به عنوان وزن بیشتری (سه برابر) داده ایم. به همین جهت score برای نتیجه ی اول بیشتر از دومی است.

```
POST /movies/_search
{
  "query": {
    "bool": {
      "should": [
        {
          "match": {
            "genres": "drama"
          }
        },
        {
          "match": {
            "title": {
              "query": "kill",
              "boost": 3
            }
          }
        }
      ]
    }
  }
}
```

```
3   "took" : 0,
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" : 5,
14      "relation" : "eq"
15    },
16    "max_score" : 4.035888,
17    "hits" : [
18      {
19        "_index" : "movies",
20        "_type" : "movie",
21        "_id" : "3",
22        "_score" : 4.035888,
23        "_source" : {
24          "title" : "To Kill a Mockingbird",
25          "director" : "Robert Mulligan",
26          "year" : 1962,
27          "genres" : [
28            "Crime",
29            "Drama",
30            "Mystery"
31          ]
32        }
33      },
34      {
35        "_index" : "movies",
36        "_type" : "movie",
37        "_id" : "5",
38        "_score" : 3.8003933,
39        "_source" : {
40          "title" : "Kill Bill: Vol. 1",
41          "director" : "Quentin Tarantino",
42          "year" : 2003,
43          "genres" : [
44            "Action",
45            "Crime",
46            "Drama"
47          ]
48        }
49      }
50    ]
51  }
52 }
```

در کوئری فوق، چند تا میچینگ را با هم ترکیب کرده ایم. عبارت should به معنی Or شدن حاصل دو میچینگ است که میچ شدن عنوان ضریب ۳ دارد. ۵ تا فیلم به عنوان نتیجه آمده اند که اولی به علت داشتن هر دو مخصوصا میچ در عنوان، بیشترین امتیاز را دریافت کرده است.

درون یک عبارت bool هر کدام از عبارت های must, must not و should تنها یک بار باید ظاهر شود. (جمله مهم دستور کار)

```
GET /movies/_search
{
  "query": {
    "term": {
      "year": {
        "value": "2003",
        "boost": 1.0
      }
    }
  }
}
```

```
1 #! Elasticsearch built-in security features are not enabled. Without authentication, your
2 cluster could be accessible to anyone. See https://www.elastic.co/guide/en/elasticsearch
3 /reference/7.16/security-minimal-setup.html to enable security.
4
5 {
6   "took": 0,
7   "timed_out": false,
8   "_shards": {
9     "total": 1,
10    "successful": 1,
11    "skipped": 0,
12    "failed": 0
13  },
14  "hits": {
15    "total": {
16      "value": 1,
17      "relation": "eq"
18    },
19    "max_score": 1.0,
20    "hits": [
21      {
22        "_index": "movies",
23        "_type": "movie",
24        "_id": "5",
25        "_score": 1.0,
26        "_source": {
27          "title": "Kill Bill: Vol. 1",
28          "director": "Quentin Tarantino",
29          "year": 2003,
30          "genres": [
31            "Action",
32            "Crime",
33            "Thriller"
34          ]
35        }
36      }
37    ]
38  }
39 }
```

برای سرچ غیر متنی مثلا فیلتر تصویر بالا روی سال ساخت، از term استفاده می کنیم. برای گشتن به دنبال مقدار دقیق برای فیلدهای متنی نیز می توان از term استفاده کرد. یعنی در حالتی که میچ شدن به قسمتی از رشته کافی نیست و کل رشته باید برابری داشته باشد.

```
PUT movies/_doc/7
{
  "title": "Kill your Darlings",
  "director": "John Krokidas",
  "year": 2013,
  "genres": ["Romance", "Drama"]
}
```

```
GET /movies/_search
{
  "query": {
    "bool": {
      "must": [
        {
          "match": {
            "title": "kill"
          }
        }
      ],
      "filter": {
        "bool": {
          "must": [
            {
              "range": {
                "year": {
                  "gte": 1960
                }
              }
            },
            {
              "term": {
                "genres": {
                  "value": "drama"
                }
              }
            }
          ]
        }
      }
    }
  }
}
```

```
2 {
3   "took": 532,
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": 2,
14      "relation": "eq"
15    },
16    "max_score": 1.1120702,
17    "hits": [
18      {
19        "_index": "movies",
20        "_type": "movie",
21        "_id": "7",
22        "_score": 1.1120702,
23        "_source": {
24          "title": "Kill your Darlings",
25          "director": "John Krokidas",
26          "year": 2013,
27          "genres": [
28            "Romance",
29            "Drama"
30          ]
31        }
32      },
33      {
34        "_index": "movies",
35        "_type": "movie",
36        "_id": "3",
37        "_score": 1.0096802,
38        "_source": {
39          "title": "To Kill a Mockingbird",
40          "director": "Robert Mulligan",
41          "year": 1962,
42          "genres": [
43            "Drama"
44          ]
45        }
46      }
47    ]
48  }
49 }
```

در کوئری فوق، فیلترینگ به همراه مچینگ اعمال کرده ایم به این صورت که داشتن ژانر درام در ژانر ها و سال ساخت بعد از ۱۹۶۰ را به عنوان فیلتر ها در نظر گرفته ایم و سپس مچ شدن عنوان با کلمه Kill را اعمال می کنیم. در مجموع کوئری ها به دو بخش متنی (matching) و غیر متنی (filtering) تقسیم می شوند که می توان این دو را ترکیب هم کرد یا حتی اصلا یکی را اعمال نکرد. مثلاً با match_all می توانیم سرچ متنی انجام ندهیم.

در این قسمت، برای داشتن حداقل دو نتیجه، فیلم Kill Your Darlings به رکوردها اضافه شد و در نتیجه نیز دیده می شود.

```
PUT movies/_doc/8
{
  "title": "Kill Bill",
  "director": "Quentin Tarantino",
  "year": 2003,
  "genres": ["Action", "Crime", "Thriller"]
}

GET /movies/_search
{
  "query": {
    "match_phrase": {
      "title": {
        "query": "Kill Bill",
        "slop": 2
      }
    }
  }
}
```

```
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" : 2,
14     "relation" : "eq"
15   },
16   "max_score" : 2.8190997,
17   "hits" : [
18     {
19       "_index" : "movies",
20       "_type" : "movie",
21       "_id" : "8",
22       "_score" : 2.8190997,
23       "_source" : {
24         "title" : "Kill Bill",
25         "director" : "Quentin Tarantino",
26         "year" : 2003,
27         "genres" : [
28           "Action",
29           "Crime",
30           "Thriller"
31         ]
32       }
33     },
34     {
35       "_index" : "movies",
36       "_type" : "movie",
37       "_id" : "5",
38       "_score" : 2.2779927,
39       "_source" : {
40         "title" : "Kill Bill: Vol. 1",
41         "director" : "Quentin Tarantino",
42         "year" : 2003,
43         "genres" : [
44           "Action",
```

در کوئری فوق، نیاز به سرچ متنی دقیق را برطرف کرده ایم که با استفاده از match_phrase گفته ایم که مقادیری انتخاب شوند که عنوانشان دقیقاً عبارت Kill Bill را دارد ولی می توانند نهایتاً ۳ تا space نیز داشته باشند. (slop=2)

برای این مرحله هم یک رکورد با عنوان Kill Bill با ۳ تا Space اضافه شده که میبینیم در نتایج نیز نمایش داده می شود.

```
GET movies/_search
{
  "query": {
    "term": {
      "director.keyword": "Francis Ford Coppola"
    }
  }
}

GET /movies/_search
{
  "query": {
    "term": {
      "director": {
        "value": "Francis Ford Coppola"
      }
    }
  }
}
```

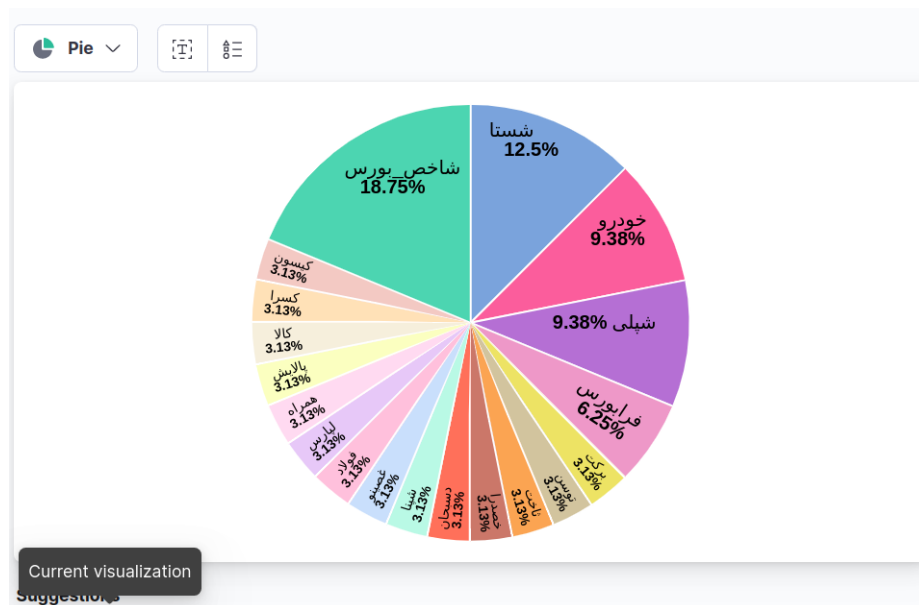
```
1 #! Elasticsearch built-in security features are not enabled. Without authentication, your
2 cluster could be accessible to anyone. See https://www.elastic.co/guide/en/elasticsearch
3 /reference/7.16/security-minimal-setup.html to enable security.
4
5 {
6   "took": 0,
7   "timed_out": false,
8   "shards": {
9     "total": 1,
10    "successful": 1,
11    "skipped": 0,
12    "failed": 0
13  },
14  "hits": {
15    "total": {
16      "value": 2,
17      "relation": "eq"
18    },
19    "max_score": 1.3862942,
20    "hits": [
21      {
22        "_index": "movies",
23        "_type": "movie",
24        "_id": "1",
25        "_score": 1.3862942,
26        "_source": {
27          "title": "The Godfather",
28          "director": "Francis Ford Coppola",
29          "year": 1972
30        }
31      },
32      {
33        "_index": "movies",
34        "_type": "movie",
35        "_id": "4",
36        "_score": 1.3862942,
37        "_source": {
38          "title": "Apocalypse Now",
39          "director": "Francis Ford Coppola",
40          "year": 1979,
41          "genres": [
42            "Drama",
```

مشکلی که وجود دارد این است که با اجرای کوئری دوم در تصویر بالا نتیجه ای نمایش داده نمی شود. آن هم به این دلیل است که الاستیک به طور خودکار فیلد های متنی را پردازش می کند و به کلمات و بخش های تشکیل دهنده شان تقسیمشان می کند. به همین علت، باید مشخص کنیم که دقیقاً به دنبال همین عبارت هستیم و این کار را با جایگزین کردن کوئری دوم با کوئری اول انجام می دهیم. یعنی می گوئیم که مقدار Francis Ford Coppola یک نوع keyword در جستجوی ماست و دقیقاً به دنبال آن عبارت هستیم. نتیجه، دیدن هر دو فیلم به کارگردانی این فرد در حاصل کوئری است.

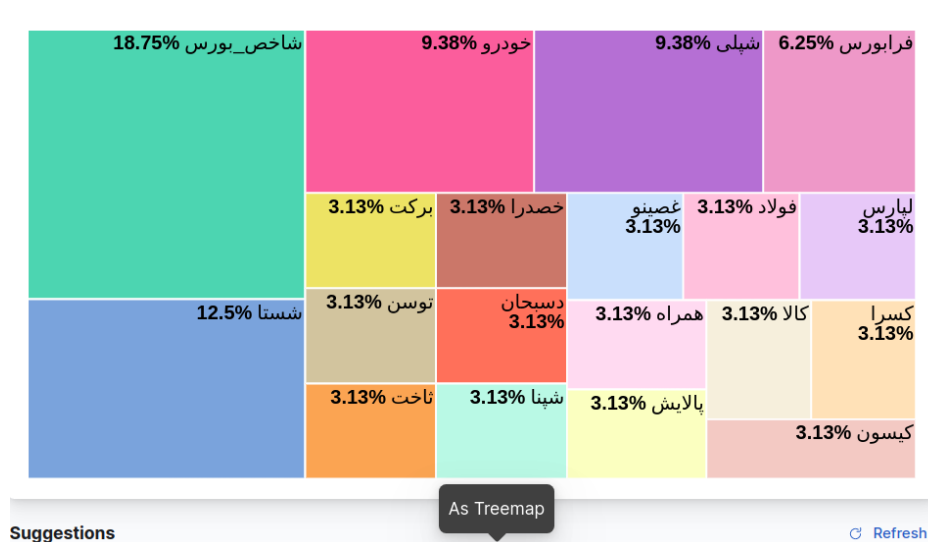
بخش دوم) جمع آوری توییت ها:

در این بخش، فایل توییت های قرار داده شده در گروه را وارد الستیک سرچ کردیم. بدین وسیله دیتای ۵۰۰ توییت وارد شد. سپس با ساختن Index pattern برای داده ورودی و ساختن یک dashboard برای visualization، داده ها و هشتگ های موجود در توییت ها و نحوه توزیع آن ها به دو روش Tree Chart (تقریباً مشابه تابلو های بورسی) و Pie Chart بصری سازی شد.

- Pie Chart:



- Tree Chart:



- احتمالا به علت ورژن متفاوت کیبانا، برای من امکان رسم Tag Cloud نبود و کیبانا آپشن رسم آن را در اختیار من نمی گذاشت. به همین علت آن را با Tree Chart جایگزین کردم.
- متاسفانه به علت نوشتن گزارش در Google Docs امکان رعایت چارچوب گزارش وجود نداشت. از این بابت عذرخواهی می کنم.