

Elastic

ارائه دهندگان

محمد خویی | علی مرتضوی منش

استاد مربوطه | دکتر سید مهدی وحیدی پور

بهار 97





ElasticSearch یک موتور جستجو و تحلیل منعطف، قدرتمند، متن باز، توزیع شده، دسترسی بالا و بلادرنگ می‌باشد که هسته شاخص‌گذار آن کتابخانه **Lucene** می‌باشد.

الاستیک سرچ یک موتور جستجوی متن با امکانات فراوان از قبیل راحتی استفاده، مقیاس پذیری بالا، فیلترهای متنوع جستجو، تحلیلگرهای اختصاصی و سفارشی برای متون فارسی قبل از ایندکس و ذخیره در دیتابیس، دسته بندی نتایج و امکان گروه بندی و فیلتر ثانویه روی آنها، جستجوی فازی و تقریبی و مانند آنها را دارد.

از ابتدا به منظور استفاده در محیط‌های توزیع شده پیاده سازی شده است، جایی که اتکاپذیری و مقیاس پذیری باید وجود داشته باشد، ES توانایی حرکت آسان ماوراء جستجوی کاملاً متنی ساده را به شما می‌دهد و رابطه‌ی اکثر زبان‌های برنامه نویسی معروف را ارائه می‌دهد، ES وعده‌های بی‌حد و حصر استفاده از فناوری جستجو را ارائه می‌کند.

قابلیت‌هایی دیگر

چندین شاخصی

مبتنی بر سند

مدیریت مغایرت

ماندگاری در هر عملیات

عدم استفاده از قالب ثابت

یک خوشه می‌تواند میزبان چند شاخصی که جدا از هم و یا به صورت یک گروه می‌باشند، باشد.

قابلیت‌هایی دیگر

چندین شاخصی

مبتنی بر سند

مدیریت مغایرت

ماندگاری در هر عملیات

عدم استفاده از قالب ثابت

ذخیره سازی موجودیت‌های پیچیده دنیای واقعی به صورت اسناد **JSON** انجام می‌گیرد. تمام فیلدها به صورت پیش فرض شاخص‌گذاری می‌شوند و تمام شاخص‌ها در یک درخواست می‌توانند استفاده شوند.

قابلیت‌های دیگر

چندین شاخصی

مبتنی بر سند

مدیریت مغایرت

ماندگاری در هر عملیات

عدم استفاده از قالب ثابت

یک کنترل کننده نسخه به منظور ممانعت از نابودی
اطلاعات در زمان تغییرات همزمان

قابلیت‌هایی دیگر

چندین شاخصی

مبتنی بر سند

مدیریت مغایرت

ماندگاری در هر عملیات

عدم استفاده از قالب ثابت

برای ES در ابتدا **امنیت** داده مهم است. تغییرات اسناد در گزارش‌های تراکنش‌ها در چندین نود در خوشه مورد نظر ثبت می‌شود تا امکان از دست رفتن داده به حداقل برسد.

رابط RESTful : استفاده از رابط RESTful به گونه‌ای که از JSON بر روی پروتکل HTTP استفاده می‌کند.

قابلیت‌هایی دیگر

چندین شاخصی

مبتنی بر سند

مدیریت مغایرت

ماندگاری در هر عملیات

عدم استفاده از قالب ثابت

ES به صورت اتوماتیک ساختار داده مورد نظر را از اسناد JSON استخراج کرده و پس از شاخص‌گذاری آن را قابل جستجو قرار می‌دهد.

سپس بعدها با اعمال دانش خاص منظوره داده‌های شما، نحوه شاخص‌گذاری داده‌های شما را تعیین می‌کند.

Lucene



Lucene این پروژه در سال ۱۹۹۹ نوشته شد و در سایت source forge به صورت متن باز ارائه شد و بعدها به بنیاد آپاچی ملحق شد. این نرم افزار یک کتابخانه نرم افزاری بازیابی اطلاعات می باشد. این نرم افزار توسط **java** نوشته شده است. این نرم افزار به زبان های برنامه نویسی دیگر از جمله دلفی، پرل، # C، C ++، پایتون، روبي و PHP پورت شده است. این نرم افزار برای هر نرم افزاری که نیاز به شاخص گذاری و جستجوی متن به صورت کامل داشته باشد بسیار مفید خواهد بود. **Lucene** به صورت گسترده ای مورد قبول واقع شده و در موتور جستجوهای اینترنتی و محلی و جستجوی تک سایت مورد استفاده قرار می گیرد.

Lucene



متن در معماری منطقی هسته Lucene ایده ای نهفته است . این ایده متذکر می شود که یک سند حاوی فیلدهایی از متن هستند . این انعطاف پذیری به API های این نرم افزار اجازه می دهد که مستقل از فرمت فایل باشند . متن از هر فرمت Open Document،HTML،WORD و ... می تواند استخراج شده و شاخص گذاری شود . این عمل بر روی تصاویر غیر قابل ممکن است .

The Lucene logo is written in a stylized, cursive font with a yellow-to-orange gradient and a drop shadow effect.

Lucene شامل توابع و کتابخانه هایی برای خزش و پارس نمودن فایل های Html می باشد ولی پروژه های دیگری این قابلیت ها را به ارائه نموده اند. اگر به لیست بانکهای اطلاعاتی برتر دنیا هم در سایت db-engines نگاهی بیندازید، این بانک را جزء ده بانک اطلاعاتی مطرح امروزمین خواهید یافت.



آشنایی با مفاهیم الاستیک سرچ

همانطور که گفتیم الاستیک سرچ یک موتور منبع باز ، جستجوی کامل متن و تجزیه و تحلیل بسیار مقیاس پذیر است. و به شما امکان ذخیره، جستجو و آنالیز حجم عظیمی از داده‌ها را در زمان اندکی می‌دهد. در ادامه چند اصطلاح پرکاربرد در ES را معرفی و شرح می‌دهیم که برای شروع کار با این موتور الزامی است :



آشنایی با مفاهیم الاستیک سرچ

Near Real Time

ES یک پلتفرم جستجوی نزدیک به زمان واقعی است این بدان معنی است که مدت زمانی که طول می کشد که یک سند از مرحله ایندکس به مرحله جستجو برسد بسیار ناچیز (Refresh Time) است.



آشنایی با مفاهیم الاستیک سرچ

Cluster

یک کلاستر مجموعه‌ای از یک یا تعداد بیشتری **نود** (Server) می‌باشد که این نودها در کنار یکدیگر **کل داده** ها را نگه داری می‌کنند و امکان **جستجو** و **ایندکس** را روی تمام نودها فراهم می‌کنند. یک کلاستر **توسط نام خود** **شناسایی** می‌شود و در حالت پیش فرض در ES نام کلاستر ElasticSearch می‌باشد. این نام بسیار **مهم** است زیرا **یک نود تنها با داشتن نام کلاستر مورد نظر میتواند به آن متصل شود** و روشن است که روی یک **بستر شبکه**، نام کلاسترها باید **متمايز** از دیگری باشد.



آشنایی با مفاهیم الاستیک سرچ

Node

هر نود می‌تواند یک سیستم جدا باشد که قسمتی از یک کلاستر را تشکیل می‌دهد. نودها محل ذخیره داده‌ها هستند و عملیات Search و Index بروی آنها انجام می‌شود. نودها هم مانند کلاستر توسط نام خود شناسایی می‌شوند و در حالت پیش فرض اگر نام نود توسط کاربر تعیین نشود ES از یک نام به صورت تصادفی از بین اسامی افراد و شخصیت‌های مشهور به آن نسبت می‌دهد و نام هر نود در شناسایی آن ضروری است و در حالت پیش‌فرض تمام نودهای ایجاد شده به کلاستر پیش فرض ES یعنی ElasticSearch اضافه می‌شوند.



آشنایی با مفاهیم الاستیک سرچ

Index

هر ایندکس، مجموعه‌ای از اسناد است که دارای ویژگی‌های تقریباً مشابه هستند و می‌توانیم روی یک کلاستر تنها به تعداد دلخواه ایندکس داشته باشیم و هر ایندکس توسط نام خود شناسایی و تعیین می‌شود و برای انجام عملیات ایندکس، جستجو، حذف و بروزرسانی دانستن نام ایندکس ضروری است.



آشنایی با مفاهیم الاستیک سرچ

Document

یک سند واحد پایه اطلاعات است که می‌تواند ایندکس شود و اسناد در ES قابل تبدیل شدن به فرمت JSON هستند.



آشنایی با مفاهیم الاستیک سرچ

Shard & Replicas

یک ایندکس می‌تواند **حجم عظیمی از داده** را ذخیره کند که از محدودیت‌های سخت افزاری یک سرور تجاوز می‌کند. بطور مثال یک ایندکس شامل یک میلیارد سند، به فضایی در حدود 1TB جهت ذخیره سازی نیاز دارد که ممکن است **بیشتر از ظرفیت یک سرور** باشد و یا اگر هم امکان ذخیره سازی این حجم از **داده** را داشته باشد **مدت زمان پاسخ** به درخواست‌های جستجو را به **مراتب افزایش** و در نتیجه باعث کاهش **کارایی کلاستر** می‌شود. لذا جهت رفع این مشکل ES توانایی **تقسیم** ایندکس‌ها به **چندین قسمت** را دارد که هر قسمت یا **تکه Shard** نامیده می‌شود و در زمان تعریف ایندکس براحتی می‌توان تعداد shard های مورد نیاز را تعیین کرد و هر **Shard** کاملاً **کاربردی و مستقل** عمل می‌کند و می‌تواند روی هر نود کلاستر قرار بگیرد.



آشنایی با مفاهیم الاستیک سرچ

Sharding به دو دلیل اصلی مهم است:

حجم داده‌های روی Shardها کم می‌شود و داده‌ها بین Shardها تقسیم می‌شوند.

امکان اجرای عملیات بصورت توزیع شده روی تمامی نودها فراهم می‌شود که در نتیجه باعث افزایش عملکرد/توان می‌شود.



آشنایی با مفاهیم الاستیک سرچ

دلایل مهم استفاده از Replication

در محیط‌هایی مانند بستر شبکه که در آن امکان شکست و Fail شدن سیستم در هر زمان را می‌توان انتظار داشت بسیار توصیه می‌شود که یک مکانیزم Fialover جهت حفظ داده‌ها و کارایی سیستم وجود داشته باشد. به همین دلیل ES به شما اجازه می‌دهد تا یک یا چند کپی از هر Shard را روی یک سیستم دیگر به نام Replica Shard داشته باشیم.

دسترس پذیری روی Shard های Fail شده



آشنایی با مفاهیم الاستیک سرچ

در خلاصه این بخش:

هر ایندکس می‌تواند داده‌های خود را روی چند Shard ذخیره کند. یک ایندکس می‌تواند چند و یا صفر Shard کپی داشته باشد. تعداد Shard ها و Replica ها می‌تواند در لحظه تعریف ایندکس تعیین شود و بعد از ایجاد ایندکس، فقط تعداد Replica ها قابل تغییر است. در ES و بصورت پیش فرض هر ایندکس روی ۵ Shard اصلی ایجاد می‌شود و هر Shard اصلی یک Replica دارد.

به طور مثال اگر کلاستر مورد نظر ما ۲ نود داشته باشد، ایندکس شما شامل ۵ Shard اصلی و ۵ Shard Replica است.

انواع رابط (API) های rest full

رابط شاخص

رابط سند

رابط جستجو

رابط کلاستر

رابطه (API) شاخص

ایجاد یک شاخص به همراه مشخص کردن ساختار داده ای سند های آن و تنظیمات مربوطه آن شاخص

مثال

(ایجاد یک شاخص)

```
1  PUT /my_index_name
2  {
3      "settings": {
4          "number_of_replicas": 1,
5          "number_of_shards": 3,
6          "analysis": {},
7          "refresh_interval": "1s"
8      },
9      "mappings": {
10         "my_type_name": {
11             "properties": {
12                 "title": {
13                     "type": "text",
14                     "analyzer": "english"
15                 }
16             }
17         }
18     }
19 }
```

رابط (API)
شاخص

تغییر یک شاخص، بخش تنظیمات

مثال تغییریک

شاخص

(بخش تنظیمات)

```
1  PUT /my_index_name/_settings
2  {
3      "index": {
4          "refresh_interval": "-1",
5          "number_of_replicas": 0
6      }
7  }
```

رابط (API)
شاخص

تغییر یک شاخص، بخش ساختار داده ای

مثال تغییریک شاخص

(ساختار داده ای)

```
1  PUT /my_index_name/_mapping/my_type_name
2  {
3      "my_type_name": {
4          "properties": {
5              "tag": {
6                  "type": "keyword"
7              }
8          }
9  }
```

رابط (API)
شاخص

دریافت بخش تنظیمات
مربوط به یک شاخص:

```
1 GET /my_index_name/_settings
```

دریافت بخش ساختار داده ای
مربوط به یک شاخص:

رابط (API)
شاخص

```
1 GET /my_index_name/_mapping
```

حذف یک شاخص:

رابط (API)
شاخص

```
1 DELETE /my_index_name
```

رابط (API)
شاخص

باز کردن و بستن یک شاخص جهت
مدیریت مصرف حافظه و پردازنده:

- 1 POST /my_index_name/_close
- 2 POST /my_index_name/_open

**رابطہ (API)
سند**

ایجاد یک سند با قابلیت ساخت کلید
اصلی خودکار

مثال

(ایجاد یک سند)

```
1  POST /my_index_name/my_type_name
2  {
3  "title":
4  |   "Elastic is funny",
5  |   "tag": [ "lucene" ]
6  }
```

رابطہ (API) سند

ایجاد و تغیر یک سند مشخص بر
اساس یک کلید مشخص :

مثال

(ایجاد و تغییر یک سند)

```
1  PUT /my_index_name/my_type_name/12abc
2  {
3  |   "title":
4  |       "Elastic is funny",
5  |   "tag": [ "lucene" ]
6  }
```

رابطہ (API) سند

حذف یک سند مشخص بر اساس یک کلید
مشخص:

```
1 DELETE /my_index_name/my_type_name/12abc
```

رابطه (API) جستجو

جستجو بر روی تمامی شاخص ها و
تمامی انواع سند موجود در آنها:

```
1 GET /_search
2 {
3   | query{}
4 }
```

رابطه (API)
جستجو

جستجو بر روی یک شاخص خاص با
مشخص کردن نوع اسناد و یا مشخص
نکردن نوع اسناد (جستو جو بر روی
تمام انواع موجود سند در شاخص):

مثال

```
1  GET /my_index_name/_search
2  = {
3      query{}
4  }
5  }
6  GET /my_index_name/type_name_1,type_name_2
7  = {
8      query{}
9  }
10 }
```

رابط (API)
جستجو

نمونه از رابط search


```
1 GET /_search
2 {
3   "query": {
4     "bool": {
5       "must": [
6         {
7           "match": {
8             "title": "smith"
9           }
10        }
11      ],
12      "filter": [
13        {
14          "exists": {
15            "field": "title"
16          }
17        }
18      ]
19    }
20  },
```

```
21     "size": 20,
22     "from": 100,
23     "_source": [
24       "title",
25       "id"
26     ],
27     "sort": [
28       {
29         "_id": {
30           "order": "desc"
31         }
32       }
33     ]
34   }
```



فیلم‌های آموزشی

جهت دریافت فیلم‌های آموزشی به کانال تلگرامی

[@elastic_learn](https://t.me/elastic_learn)

مراجعه کنید.

**برخی موارد
پیشرفته تر**

تعریف :

تعریف سایت رهنما:

Analysis is the process of converting text, like the body of any email, into *tokens* or *terms* which are added to the inverted index for searching. Analysis is performed by an analyzer which can be either a built-in analyzer or a custom analyzer defined per index.

مشاهده ی نحوه ی عملکرد یکی از تحلیلگر های از پیش تعریف شده

"The QUICK brown foxes jumped over the lazy dog!"

[quick, brown, fox, jump, over, lazy, dog]

فیلتر گذاری کارکترها (character filter)

متن پایه و خام را دریافت میکند و با تغییر و یا حذف و یا اضافه کردن کارکترها متن پایه برای فرایند تشخیص کلمات آماده می کند .

تعریف سایت راهنما:

A character filter receives the original text as a stream of characters and can transform the stream by adding, removing, or changing characters. For instance, a character filter could be used to convert Hindu-Arabic numerals (۰۱۲۳۴۵۶۷۸۹) into their Arabic-Latin equivalents (0123456789), or to strip HTML elements like from the stream.

An analyzer may have zero or more [character filters](#), which are applied in order.

تشخیص و استخراج کلمات (tokenizer)

کاراکترهای ورودی را دریافت میکند و با تشخیص کلمه، کلمه شناخته شده را جدا سازی می کند

تعریف سایت راهنما:

A tokenizer receives a stream of characters, breaks it up into individual tokens (usually individual words), and outputs a stream of tokens. For instance, a [whitespace](#) tokenizer breaks text into tokens whenever it sees any whitespace. It would convert the text "Quick brown fox!" into the terms [Quick, brown, fox!].

The tokenizer is also responsible for recording the order or position of each term and the start and end character offsets of the original word which the term represents.

An analyzer must have exactly one [tokenizer](#).

فیلتر گذاری بر روی کلمات (token filter)

کلمات را دریافت میکند و فیلترهایی را بر روی آن اعمال میکند که میتواند موجب تغییر حذف و یا ایجاد کلمات شود.

تعریف سایت راهنما:

A token filter receives the token stream and may add, remove, or change tokens. For example, a [lowercase](#) token filter converts all tokens to lowercase, a [stop](#) token filter removes common words (stop words) like the from the token stream, and a [synonym](#) token filter introduces synonyms into the token stream.

Token filters are not allowed to change the position or character offsets of each token.

An analyzer may have zero or more [token filters](#), which are applied in order.

Standard Analyzer

The standard analyzer divides text into terms on word boundaries, as defined by the Unicode Text Segmentation algorithm. It removes most punctuation, lowercases terms, and supports removing stop words.

Simple Analyzer

The simple analyzer divides text into terms whenever it encounters a character which is not a letter. It lowercases all terms.

Whitespace Analyzer

The whitespace analyzer divides text into terms whenever it encounters any whitespace character. It does not lowercase terms.

```
1 PUT my_index
2 {
3   "settings": {
4     "analysis": {
5       "analyzer": {
6         "std_english": {
7           "type": "standard",
8           "stopwords": "_english_"
9         }
10      }
11    }
12  },
13  "mappings": {
14    "_doc": {
15      "properties": {
16        "my_text": {
17          "type": "text",
18          "analyzer": "standard",
19          "fields": {
20            "english": {
21              "type": "text",
22              "analyzer": "std_english"
23            }
24          }
25        }
26      }
27    }
28  }
29 }
```

تغییر یک نمونه تحلیل گر از پیش تعریف شده:

انواع حالت های انتساب یک تحلیلگر:

- انتساب به یک فیلد خاص یک سند
- انتساب به یک سند
- انتساب به یک شاخص


```
1 PUT my_index
2 {
3     "mappings": {
4         "_doc": {
5             "properties": {
6                 "title": {
7                     "type": "text",
8                     "analyzer": "standard"
9                 }
10            }
11        }
12    }
13 }
```

انتساب یک تحلیلگر (analyzer) برای یک فیلد خاص
از یک نوع سند خاص در یک شاخص:

```
1 PUT /english_example
2 {
3   "settings": {
4     "analysis": {
5       "filter": {
6         "english_stop": {
7           "type": "stop",
8           "stopwords": "_english_"
9         },
10        "english_stemmer": {
11          "type": "stemmer",
12          "language": "english"
13        }
14      },
15      "analyzer": {
16        "my_english": {
17          "char_filter": ["html_strip"],
18          "tokenizer": "standard",
19          "filter": [
20            "lowercase",
21            "english_stop",
22            "english_stemmer"
23          ]
24        }
25      }
26    }
27  }
28 }
```

ایجاد یک نمونه تحلیل گر (analyzer):

مراجع

- www.bigdata-ir.com
- www.tutorialspoint.com/elasticsearch/elasticsearch_analysis.htm
- www.tutorialspoint.com/elasticsearch/elasticsearch_document_apis.htm
- www.tutorialspoint.com/elasticsearch/elasticsearch_index_apis.htm
- www.tutorialspoint.com/elasticsearch/elasticsearch_mapping.htm
- www.tutorialspoint.com/elasticsearch/elasticsearch_search_apis.htm

باتشکر از توجه شما

