

ELK搜索高级课程

#1. 课程简介

#1.1 课程内容

ELK是包含但不限于Elasticsearch（简称es）、Logstash、Kibana 三个开源软件的组成的一个整体，分别取其首字母组成ELK。ELK是用于数据抽取（Logstash）、搜索分析（Elasticsearch）、数据展现（Kibana）的一整套解决方案，所以也称作ELK stack。

本课程分别对三个组件进行详细介绍，尤其是Elasticsearch，因为它是ELK的核心。Elasticsearch会对底层的文档、索引、搜索、聚合、集群进行介绍，从搜索和聚合分析实例来展现它的魅力。Logstash会从内部如何采集数据到指定地方来展现它数据采集的功能。Kibana则从数据绘图来展现它数据可视化的能力。

#1.2 面向人员

- java工程师：深入研究es,使得java工程师向搜索工程师迈进。
- 运维工程师：搭建整体elk集群。不需写代码，仅需配置，即可收集服务器指标、日志文件、数据库数据，并在前端华丽展现。
- 数据分析人员：不需写代码，仅需配置kibana图表，即可完成数据可视化工作，得到想要的数据图表。
- 大厂架构师：完成数据中台的搭建。对公司数据流的处理得心应手，对接本公司大数据业务。

#1.3 课程优势

- 基于最新的elk7.3版本讲解。最新api。包含sql功能。
- 理论和实际代码相辅相成。理论结合画图讲解。代码与spring boot结合。
- 包含实际运维部署理论与实践。
- ELK整体流程项目，包含数据采集。

#1.4 学习路径

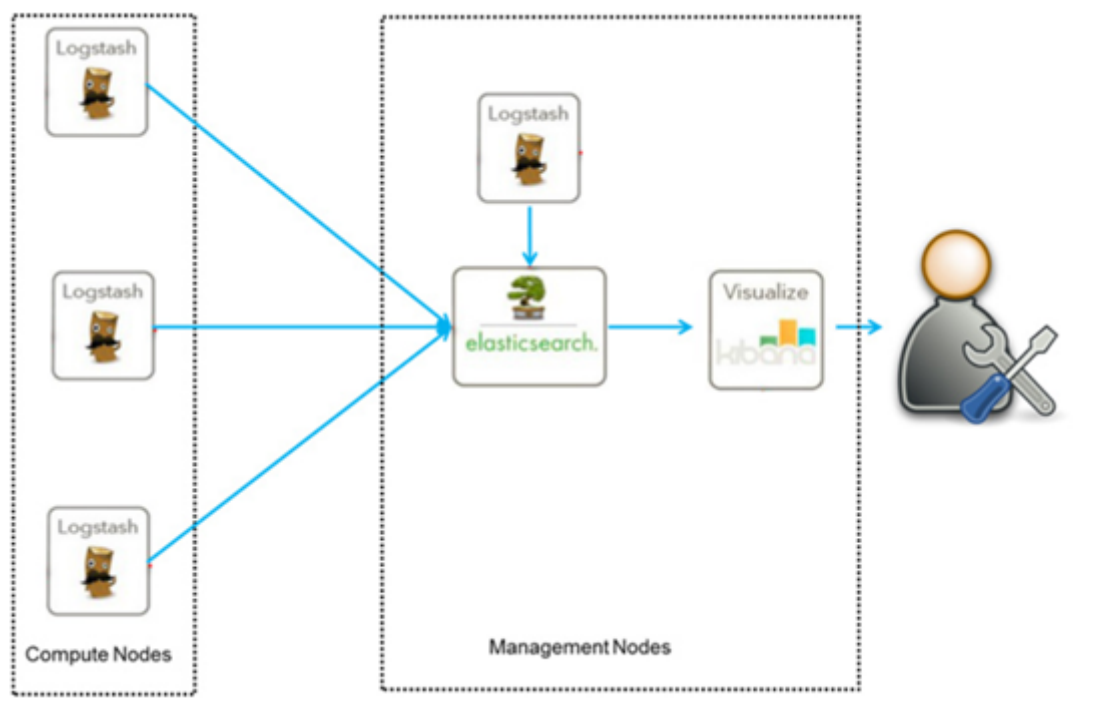
参照目录，按照介绍，es入门，文档、映射、索引、分词器、搜索、聚合。logstash、kibana。集群部署。项目实战。

每个知识点先学概念，在学rest api,最后java代码上手。

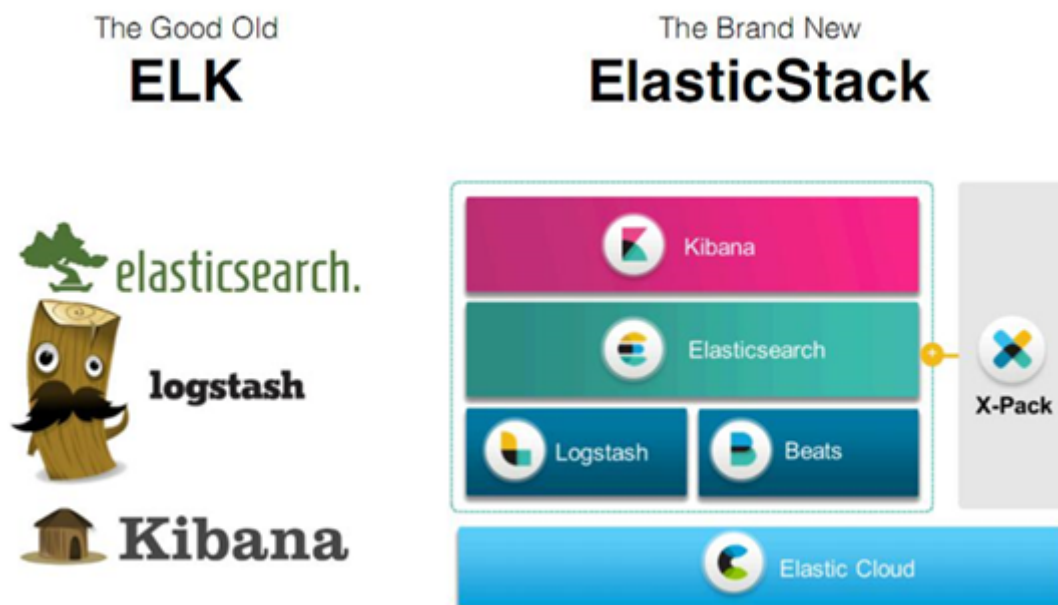
#2. Elastic Stack简介

#2.1 简介

ELK是一个免费开源的日志分析架构技术栈总称，官网 <https://www.elastic.co/cn>。包含三大基础组件，分别是Elasticsearch、Logstash、Kibana。但实际上ELK不仅仅适用于日志分析，它还可以支持其它任何数据搜索、分析和收集的场景，日志分析和收集只是更具有代表性。并非唯一性。下面是ELK架构：



随着elk的发展，又有新成员Beats、elastic cloud的加入，所以就形成了Elastic Stack。所以说，ELK是旧的称呼，Elastic Stack是新的名字。



#2.2特色

- 处理方式灵活：elasticsearch是目前最流行的准实时全文检索引擎，具有高速检索大数据的能力。
- 配置简单：安装elk的每个组件，仅需配置每个组件的一个配置文件即可。修改处不多，因为大量参数已经默认配在系统中，修改想要修改的选项即可。
- 接口简单：采用json形式RESTFUL API接受数据并响应，无关语言。
- 性能高效：elasticsearch基于优秀的全文搜索技术Lucene，采用倒排索引，可以轻易地在百亿级别数据量下，搜索出想要的内容，并且是秒级响应。

- 灵活扩展：elasticsearch和logstash都可以根据集群规模线性拓展，elasticsearch内部自动实现集群协作。
- 数据展现华丽：kibana作为前端展现工具，图表华丽，配置简单。

2.3 组件介绍

Elasticsearch

Elasticsearch 是使用java开发，基于Lucene、分布式、通过Restful方式进行交互的近实时搜索平台框架。它的特点有：分布式，零配置，自动发现，索引自动分片，索引副本机制，restful风格接口，多数据源，自动搜索负载等。

Logstash

Logstash 基于java开发，是一个数据抽取转化工具。一般工作方式为c/s架构，client端安装在需要收集信息的主机上，server端负责将收到的各节点日志进行过滤、修改等操作在一并发往elasticsearch或其他组件上去。

Kibana

Kibana 基于nodejs，也是一个开源和免费的可视化工具。Kibana可以为 Logstash 和 ElasticSearch 提供的日志分析友好的 Web 界面，可以汇总、分析和搜索重要数据日志。

Beats

Beats 平台集合了多种单一用途数据采集器。它们从成百上千或成千上万台机器和系统向 Logstash 或 Elasticsearch 发送数据。

Beats由如下组成:

Packetbeat：轻量型网络数据采集器，用于深挖网线上传输的数据，了解应用程序动态。Packetbeat 是一款轻量型网络数据包分析器，能够将数据发送至 Logstash 或 Elasticsearch。其支持ICMP (v4 and v6)、DNS、HTTP、Mysql、PostgreSQL、Redis、MongoDB、Memcache等协议。

Filebeat：轻量型日志采集器。当您要面对成百上千、甚至成千上万的服务器、虚拟机和容器生成的日志时，请告别 SSH 吧。Filebeat 将为您提供一种轻量型方法，用于转发和汇总日志与文件，让简单的事情不再繁杂。

Metricbeat：轻量型指标采集器。Metricbeat 能够以一种轻量型的方式，输送各种系统和服务统计数据，从 CPU 到内存，从 Redis 到 Nginx，不一而足。可定期获取外部系统的监控指标信息，其可以监控、收集 Apache http、HAProxy、MongoDB、MySQL、Nginx、PostgreSQL、Redis、System、Zookeeper 等服务。

Winlogbeat：轻量型 Windows 事件日志采集器。用于密切监控基于 Windows 的基础设施上发生的事件。Winlogbeat 能够以一种轻量型的方式，将 Windows 事件日志实时地流式传输至 Elasticsearch 和 Logstash。

Auditbeat：轻量型审计日志采集器。收集您 Linux 审计框架的数据，监控文件完整性。Auditbeat 实时采集这些事件，然后发送到 Elastic Stack 其他部分做进一步分析。

Heartbeat：面向运行状态监测的轻量型采集器。通过主动探测来监测服务的可用性。通过给定 URL 列表，Heartbeat 仅仅询问：网站运行正常吗？Heartbeat 会将此信息和响应时间发送至 Elastic 的其他部分，以进行进一步分析。

Functionbeat：面向云端数据的无服务器采集器。在作为一项功能部署在云服务提供商的功能即服务 (FaaS) 平台上后，Functionbeat 即能收集、传送并监测来自您的云服务的相关数据。

Elastic cloud

基于 Elasticsearch 的软件即服务(SaaS)解决方案。通过 Elastic 的官方合作伙伴使用托管的 Elasticsearch 服务。



#3. Elasticsearch是什么

#3.1搜索是什么

概念：用户输入想要的关键词，返回含有该关键词的所有信息。

场景：

1互联网搜索：谷歌、百度、各种新闻首页

2 站内搜索（垂直搜索）：企业OA查询订单、人员、部门，电商网站内部搜索商品（淘宝、京东）场景。

#3.2 数据库做搜索弊端

#3.2.1站内搜索（垂直搜索）：数据量小，简单搜索，可以使用数据库。

问题出现：

l 存储问题。电商网站商品上亿条时，涉及到单表数据过大必须拆分表，数据库磁盘占用过大必须分库（mycat）。

l 性能问题：解决上面问题后，查询“笔记本电脑”等关键词时，上亿条数据的商品名字段逐行扫描，性能跟不上。

l 不能分词。如搜索“笔记本电脑”，只能搜索完全和关键词一样的数据，那么数据量小时，搜索“笔记电脑”，“电脑”数据要不要给用户。

3.2.2 互联网搜索，肯定不会使用数据库搜索。数据量太大。PB级。

3.3 全文检索、倒排索引和Lucene

全文检索

倒排索引。数据存储时，经行分词建立term索引库。见画图。

倒排索引源于实际应用中需要根据属性的值来查找记录。这种索引表中的每一项都包括一个属性值和具有该属性值的各记录的地址。由于不是由记录来确定属性值，而是由属性值来确定记录的位置，因而称为倒排索引(inverted index)。带有倒排索引的文件我们称为倒排 **索引文件** [open in new window](#)，简称 **倒排文件** [open in new window](#) (inverted file)。

Lucene

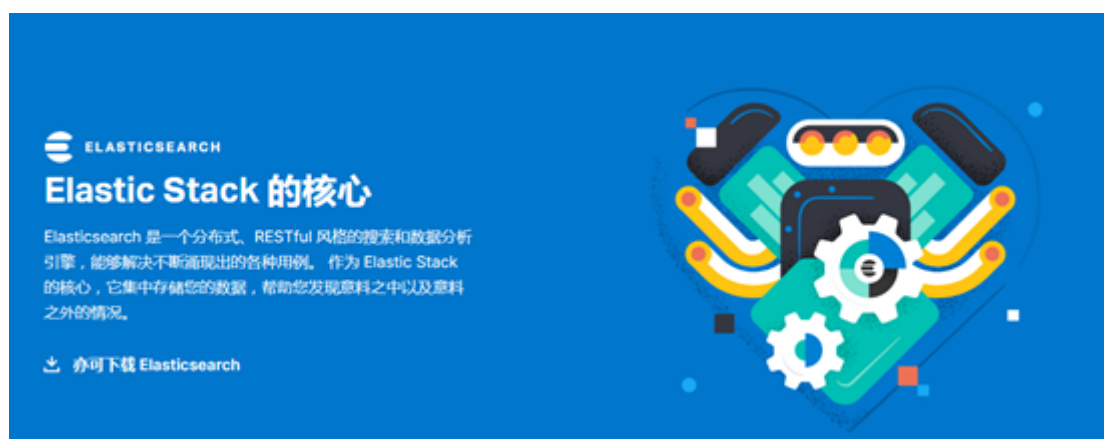
就是一个jar包，里面封装了全文检索的引擎、搜索的算法代码。开发时，引入lucene的jar包，通过api开发搜索相关业务。底层会在磁盘建立索引库。

3.4 什么是Elasticsearch

简介

ElasticSearch是一个基于Lucene的搜索服务器。它提供了一个分布式多用户能力的全文搜索引擎，基于RESTful web接口。Elasticsearch是用Java语言开发的，并作为Apache许可条款下的开放源码发布，是一种流行的企业级搜索引擎。ElasticSearch用于云计算中，能够达到实时搜索，稳定，可靠，快速，安装使用方便。官方客户端在Java、.NET (C#)、PHP、Python、Apache Groovy、Ruby和许多其他语言中都是可用的。根据DB-Engines的排名显示，Elasticsearch是最受欢迎的企业搜索引擎，其次是Apache Solr，也是基于Lucene。

官网：<https://www.elastic.co/cn/products/elasticsearch>



Elasticsearch的功能

- 分布式的搜索引擎和数据分析引擎

搜索：互联网搜索、电商网站站内搜索、OA系统查询

数据分析：电商网站查询近一周哪些品类的图书销售前十；新闻网站，最近3天阅读量最高的十个关键词，舆情分析。

- 全文检索，结构化检索，数据分析

全文检索：搜索商品名称包含java的图书select * from books where book_name like "%java%"。

结构化检索：搜索商品分类为spring的图书都有哪些，`select * from books where category_id='spring'`

数据分析：分析每一个分类下有多少种图书，`select category_id,count(*) from books group by category_id`

- 对海量数据进行近实时的处理

分布式：ES自动可以将海量数据分散到多台服务器上去存储和检索,经行并行查询，提高搜索效率。相对的，Lucene是单机应用。

近实时：数据库上亿条数据查询，搜索一次耗时几个小时，是批处理（batch-processing）。而es只需秒级即可查询海量数据，所以叫近实时。秒级。

Elasticsearch的使用场景

国外：

- 维基百科，类似百度百科，“网络七层协议”的维基百科，全文检索，高亮，搜索推荐
- Stack Overflow（国外的程序讨论论坛），相当于程序员的贴吧。遇到it问题去上面发帖，热心网友下面回帖解答。
- GitHub（开源代码管理），搜索上千亿行代码。
- 电商网站，检索商品
- 日志数据分析，logstash采集日志，ES进行复杂的数据分析（ELK技术，elasticsearch+logstash+kibana）
- 商品价格监控网站，用户设定某商品的价格阈值，当低于该阈值的时候，发送通知消息给用户，比如说订阅《java编程思想》的监控，如果价格低于27块钱，就通知我，我就去买。
- BI系统，商业智能（Business Intelligence）。大型连锁超市，分析全国网点传回的数据，分析各个商品在什么季节的销售量最好、利润最高。成本管理，店面租金、员工工资、负债等信息进行分析。从而部署下一个阶段的战略目标。

国内：

- 百度搜索，第一次查询，使用es。
- OA、ERP系统站内搜索。

Elasticsearch的特点

- 可拓展性：大型分布式集群（数百台服务器）技术，处理PB级数据，大公司可以使用。小公司数据量小，也可以部署在单机。大数据领域使用广泛。
- 技术整合：将全文检索、数据分析、分布式相关技术整合在一起：lucene（全文检索），商用的数据分析软件（BI软件），分布式数据库（mycat）
- 部署简单：开箱即用，很多默认配置不需关心，解压完成直接运行即可。拓展时，只需多部署几个实例即可，负载均衡、分片迁移集群内部自己实施。
- 接口简单：使用restful api经行交互，跨语言。
- 功能强大：Elasticsearch作为传统数据库的一个补充，提供了数据库所不能提供的很多功能，如全文检索，同义词处理，相关度排名。

可扩展性

**可以在笔记本电脑上运行。
也可以在承载了PB级数据的
成百上千台服务器上运行。**

原型环境和生产环境可无缝切换；无论 Elasticsearch 是在一个节点上运行，还是在一个包含 300 个节点的集群上运行，您都能够以相同的方式与 Elasticsearch 进行通信。

它能够水平扩展，每秒钟可处理海量事件，同时能够自动管理索引和查询在集群中的分布方式，以实现极其流畅的操作。





弹性

我们在您高飞的时候保驾护航。

相关性

搜索所有内容。找到所需的具体信息。

基于各项元素（从词频或近因到热门度等）对搜索结果进行排序。将这些内容与功能进行混合和匹配，以对向用户显示结果的方式进行微调。

而且，由于我们的大部分用户都是真实的人，Elasticsearch 具备齐全功能，可以处理包括各种复杂情况（例如拼写错误）在内的人为错误。

硬件故障、网络分割。Elasticsearch 为您检测这些故障并确保您的集群（和数据）的安全性和可用性。通过跨集群复制功能，辅助集群可以作为热备份随时投入使用。Elasticsearch 运行在一个分布式的环境中，从设计之初就考虑到了这一点，目的只有一个，让您永远高枕无忧。

3.5 elasticsearch核心概念

3.5.1 lucene和elasticsearch的关系

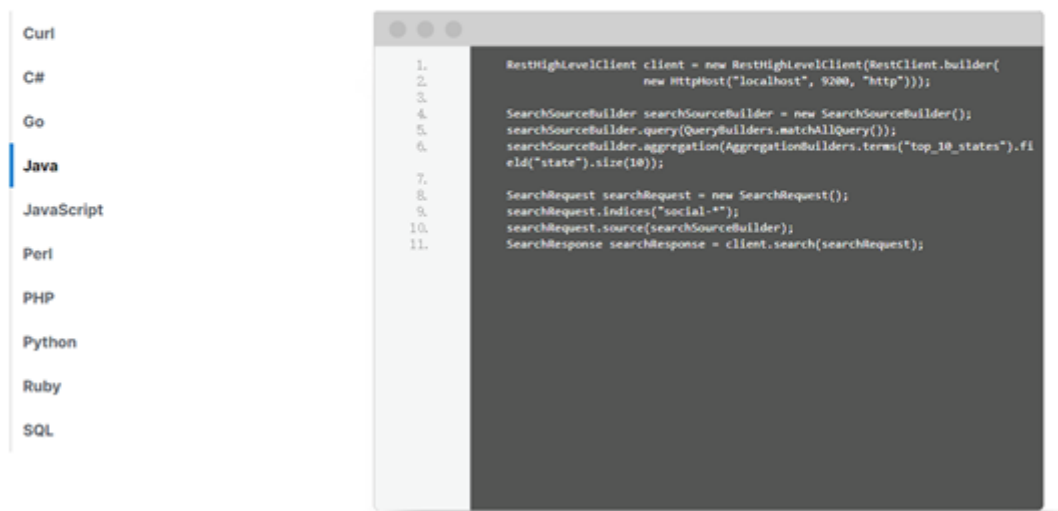
Lucene：最先进、功能最强大的搜索库，直接基于lucene开发，非常复杂，api复杂

Elasticsearch：基于lucene，封装了许多lucene底层功能，提供简单易用的restful api接口和许多语言的客户端，如java的高级客户端（[Java High Level REST Clientopen in new window](#)）和底层客户端（[Java Low Level REST Clientopen in new window](#)）

客户端库

使用您自己的编程语言与Elasticsearch 进行交互

Elasticsearch 使用的是标准的 RESTful 风格的 API 和 JSON。此外，我们还构建和维护了很多其他语言的客户端，例如 [Java](#)、[Python](#)、[.NET](#)、[SQL](#) 和 [PHP](#)。与此同时，我们的[社区](#)也贡献了很多客户端。这些客户端使用起来简单自然，而且就像 Elasticsearch 一样，不会对您的使用方式进行限制。



起源：Shay Banon。2004年失业，陪老婆去伦敦学习厨师。失业在家帮老婆写一个菜谱搜索引擎。封装了lucene的开源项目，compass。找到工作后，做分布式高性能项目，再封装compass，写出了elasticsearch，使得lucene支持分布式。现在是Elasticsearch创始人兼Elastic首席执行官。

3.5.2 elasticsearch的核心概念

1 NRT (Near Realtime) : 近实时

两方面：

- 写入数据时，过1秒才会被搜索到，因为内部在分词、录入索引。
- es搜索时：搜索和分析数据需要秒级出结果。

2 Cluster: 集群

包含一个或多个启动着es实例的机器群。通常一台机器起一个es实例。同一网络下，集群名一样的多个es实例自动组成集群，自动均衡分片等行为。默认集群名为“elasticsearch”。

3 Node: 节点

每个es实例称为一个节点。节点名自动分配，也可以手动配置。

4 Index: 索引

包含一堆有相似结构的文档数据。

索引创建规则：

- 仅限小写字母
- 不能包含\、/、*、?、"、<、>、|、#以及空格等特殊符号
- 从7.0版本开始不再包含冒号
- 不能以-、_或+开头
- 不能超过255个字节（注意它是字节，因此多字节字符将计入255个限制）

5 Document: 文档

es中的最小数据单元。一个document就像数据库中的一条记录。通常以json格式显示。多个document存储于一个索引（Index）中。

```
1  book document
2
3  {
4    "book_id": "1",
5    "book_name": "java编程思想",
6    "book_desc": "从Java的基础语法到最高级特性（深入的[面向对象]
    (https://baike.baidu.com/item/面向对象)概念、多线程、自动项目构建、单元测试和调试等），本书都能逐步指导你轻松掌握。",
7    "category_id": "2",
8    "category_name": "java"
9  }
```

6 Field: 字段

就像数据库中的列（Columns），定义每个document应该有的字段。

7 Type: 类型

每个索引里都可以有一个或多个type，type是index中的一个逻辑数据分类，一个type下的document，都有相同的field。

注意：6.0之前的版本有type（类型）概念，type相当于关系数据库的表，ES官方将在ES9.0版本中彻底删除type。本教程typy都为_doc。

#8 shard: 分片

index数据过大时，将index里面的数据，分为多个shard，分布式的存储在各个服务器上面。可以支持海量数据和高并发，提升性能和吞吐量，充分利用多台机器的cpu。

#9 replica: 副本

在分布式环境下，任何一台机器都会随时宕机，如果宕机，index的一个分片没有，导致此index不能搜索。所以，为了保证数据的安全，我们会将每个index的分片进行备份，存储在另外的机器上。保证少数机器宕机es集群仍可以搜索。

能正常提供查询和插入的分片我们叫做主分片（primary shard），其余的我们就管他们叫做备份的分片（replica shard）。

es6默认新建索引时，5分片，2副本，也就是一主一备，共10个分片。所以，es集群最小规模为两台。

#3.5.3 elasticsearch核心概念 vs. 数据库核心概念

关系型数据库（比如Mysql）	非关系型数据库（Elasticsearch）
数据库Database	索引Index
表Table	索引Index（原为Type）
数据行Row	文档Document
数据列Column	字段Field
约束 Schema	映射Mapping

#4. Elasticsearch相关软件安装

#4.1. Windows安装elasticsearch

#1、安装JDK，至少1.8.0_73以上版本，验证：java -version。

#2、下载和解压缩Elasticsearch安装包，查看目录结构。

<https://www.elastic.co/cn/downloads/elasticsearch>

bin：脚本目录，包括：启动、停止等可执行脚本

config：配置文件目录

data：索引目录，存放索引文件的地方

logs：日志目录

modules：模块目录，包括了es的功能模块

plugins：插件目录，es支持插件机制

#3、配置文件：

位置：

ES的配置文件的地址根据安装形式的不同而不同：

使用zip、tar安装，配置文件的地址在安装目录的config下。

使用RPM安装，配置文件在/etc/elasticsearch下。

使用MSI安装，配置文件的地址在安装目录的config下，并且会自动将config目录地址写入环境变量ES_PATH_CONF。

elasticsearch.yml

配置格式是YAML，可以采用如下两种方式：

方式1：层次方式

```
1 path:
2   data: /var/lib/elasticsearch
3   logs: /var/log/elasticsearch
```

方式2：属性方式

```
1 path.data: /var/lib/elasticsearch
2 path.logs: /var/log/elasticsearch
```

常用的配置项如下

```
1 cluster.name:
2   配置elasticsearch的集群名称，默认是elasticsearch。建议修改成一个有意义的名称。
3 node.name:
4   节点名，通常一台物理服务器就是一个节点，es会默认随机指定一个名字，建议指定一个有意义的名称，
   方便管理
5   一个或多个节点组成一个cluster集群，集群是一个逻辑的概念，节点是物理概念，后边章节会详细介绍。
6 path.conf:
7   设置配置文件的存储路径，tar或zip包安装默认在es根目录下的config文件夹，rpm安装默认在/etc/
   elasticsearch
8 path.data:
9   设置索引数据的存储路径，默认是es根目录下的data文件夹，可以设置多个存储路径，用逗号隔开。
10 path.logs:
11   设置日志文件的存储路径，默认是es根目录下的logs文件夹
12 path.plugins:
13   设置插件的存放路径，默认是es根目录下的plugins文件夹
14 bootstrap.memory_lock: true
15   设置为true可以锁住ES使用的内存，避免内存与swap分区交换数据。
16 network.host:
17   设置绑定主机的ip地址，设置为0.0.0.0表示绑定任何ip，允许外网访问，生产环境建议设置为具体的
   ip。
18 http.port: 9200
19   设置对外服务的http端口，默认为9200。
20 transport.tcp.port: 9300  集群结点之间通信端口
21 node.master:
22   指定该节点是否有资格被选举成为master结点，默认是true，如果原来的master宕机会重新选举新的
   master。
23 node.data:
24   指定该节点是否存储索引数据，默认为true。
25 discovery.zen.ping.unicast.hosts: ["host1:port", "host2:port", "..."]
```

```
26     设置集群中master节点的初始列表。
27     discovery.zen.ping.timeout: 3s
28     设置ES自动发现节点连接超时的时间，默认为3秒，如果网络延迟高可设置大些。
29     discovery.zen.minimum_master_nodes:
30     主结点数量的最少值，此值的公式为: (master_eligible_nodes / 2) + 1，比如：有3个符合要求的主结点，那么这里要设置为2。
31     node.max_local_storage_nodes:
32     单机允许的最大存储结点数，通常单机启动一个结点建议设置为1，开发环境如果单机启动多个节点可设置大于1。
```

jvm.options

设置最小及最大的JVM堆内存大小：

在jvm.options中设置 -Xms和-Xmx：

- 1) 两个值设置为相等
- 2) 将Xmx 设置为不超过物理内存的一半。

log4j2.properties

日志文件设置，ES使用log4j，注意日志级别的配置。

#4、启动Elasticsearch：bin\elasticsearch.bat，es的特点就是开箱即，无需配置，启动即可。

注意：es7 windows版本不支持机器学习，所以elasticsearch.yml中添加如下几个参数：

```
1     node.name: node-1
2     cluster.initial_master_nodes: [ "node-1" ]
3     xpack.ml.enabled: false
4     http.cors.enabled: true
5     http.cors.allow-origin: /*/
```

#5、检查ES是否启动成功：浏览器访问 http://localhost:9200/?Pretty

```
1     {
2         "name": "node-1",
3         "cluster_name": "elasticsearch",
4         "cluster_uuid": "HqAKQ_0tQ00m8b6qU-2Qug",
5         "version": {
6             "number": "7.3.0",
7             "build_flavor": "default",
8             "build_type": "zip",
9             "build_hash": "de777fa",
10            "build_date": "2019-07-24T18:30:11.767338Z",
11            "build_snapshot": false,
12            "lucene_version": "8.1.0",
13            "minimum_wire_compatibility_version": "6.8.0",
14            "minimum_index_compatibility_version": "6.0.0-beta1"
15        },
16        "tagline": "You Know, for Search"
17    }
```

解释：

name: node名称，取自机器的hostname

cluster_name: 集群名称（默认的集群名称就是elasticsearch）

version.number: 7.3.0, es版本号

version.lucene_version: 封装的lucene版本号

#6、浏览器访问 http://localhost:9200/_cluster/health 查询集群状态

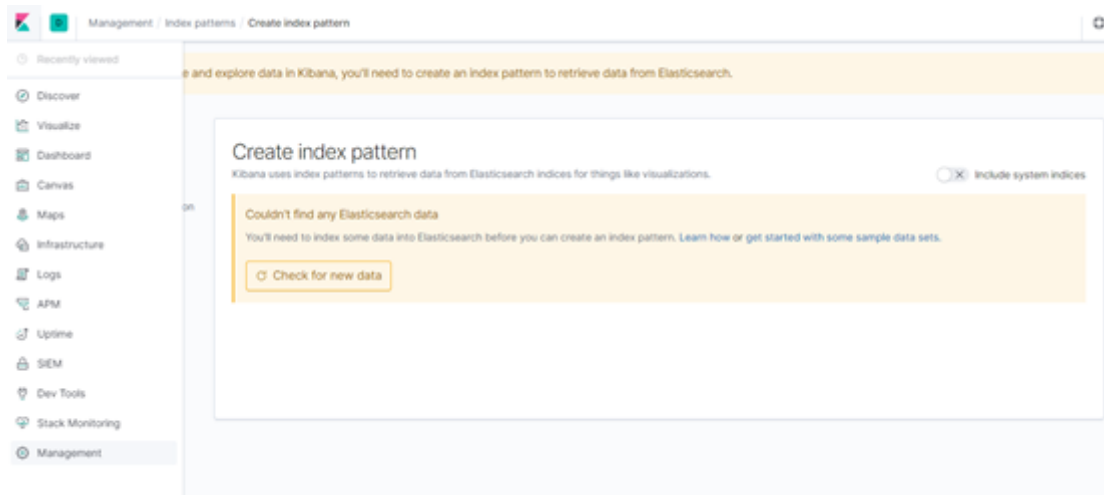
```
1  {
2      "cluster_name": "elasticsearch",
3      "status": "green",
4      "timed_out": false,
5      "number_of_nodes": 1,
6      "number_of_data_nodes": 1,
7      "active_primary_shards": 0,
8      "active_shards": 0,
9      "relocating_shards": 0,
10     "initializing_shards": 0,
11     "unassigned_shards": 0,
12     "delayed_unassigned_shards": 0,
13     "number_of_pending_tasks": 0,
14     "number_of_in_flight_fetch": 0,
15     "task_max_waiting_in_queue_millis": 0,
16     "active_shards_percent_as_number": 100
17 }
```

解释：

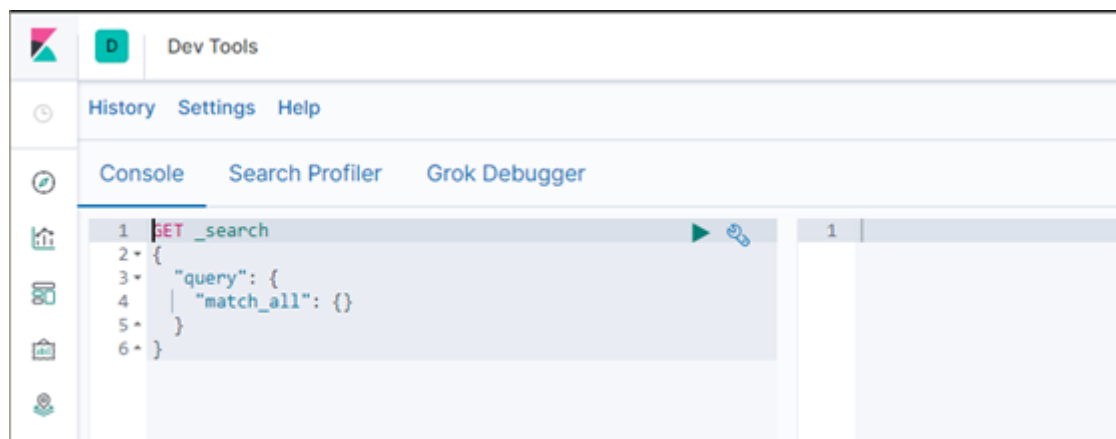
Status： 集群状态。Green 所有分片可用。Yellow所有主分片可用。Red主分片不可用， 集群不可用。

#4.2. Windows安装Kibana

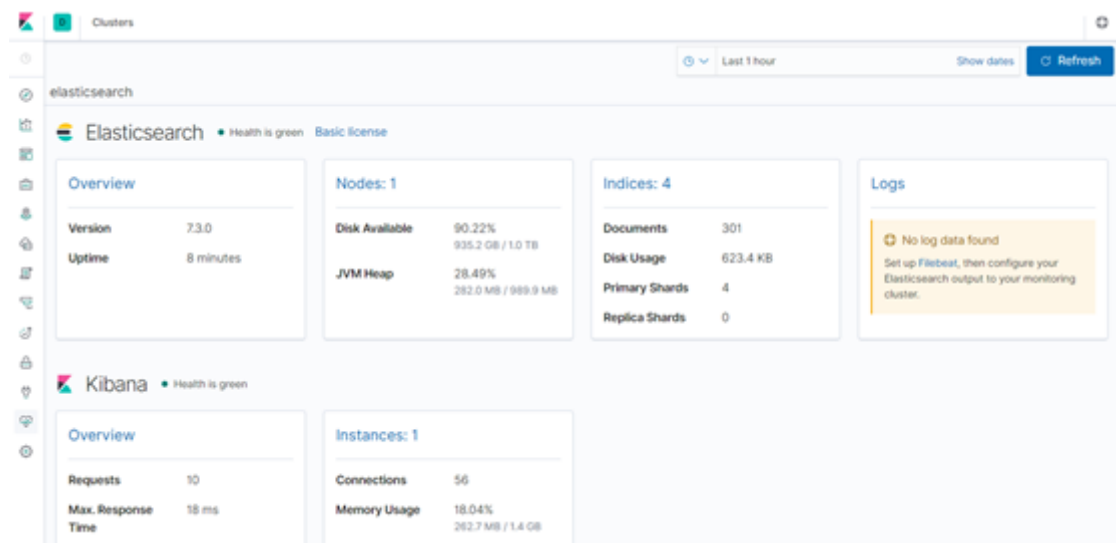
- 1、kibana是es数据的前端展现，数据分析时，可以方便地看到数据。作为开发人员，可以方便访问es。
- 2、下载，解压kibana。
- 3、启动Kibana：bin\kibana.bat
- 4、浏览器访问 <http://localhost:5601> 进入Dev Tools界面。像plsql一样支持代码提示。
- 5、发送get请求，查看集群状态GET _cluster/health。相当于浏览器访问。



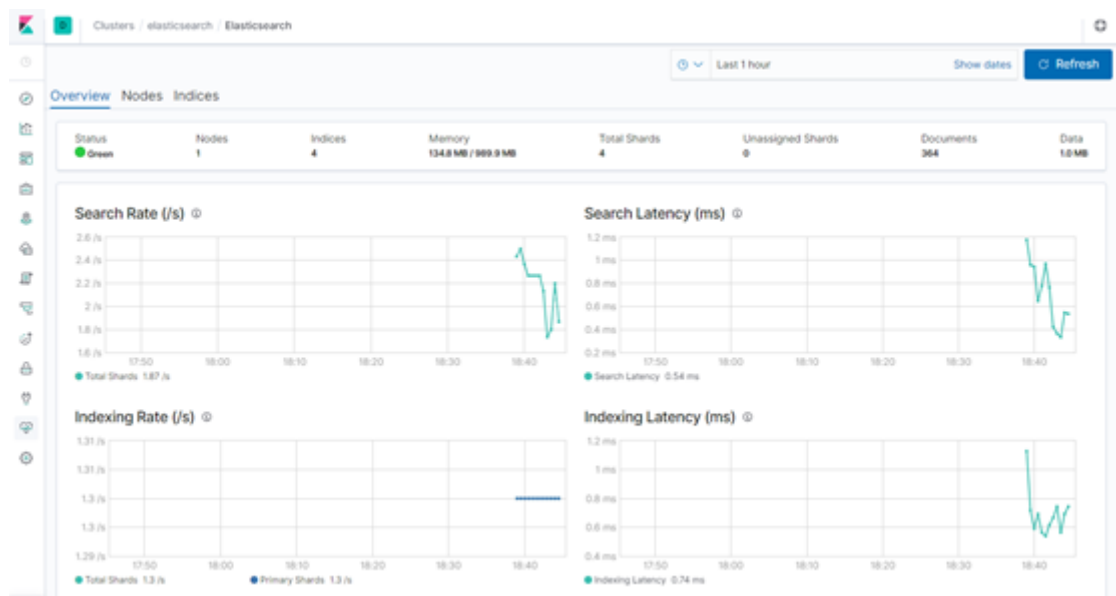
总览



Dev Tools界面



监控集群界面



集群状态（搜索速率、索引速率等）

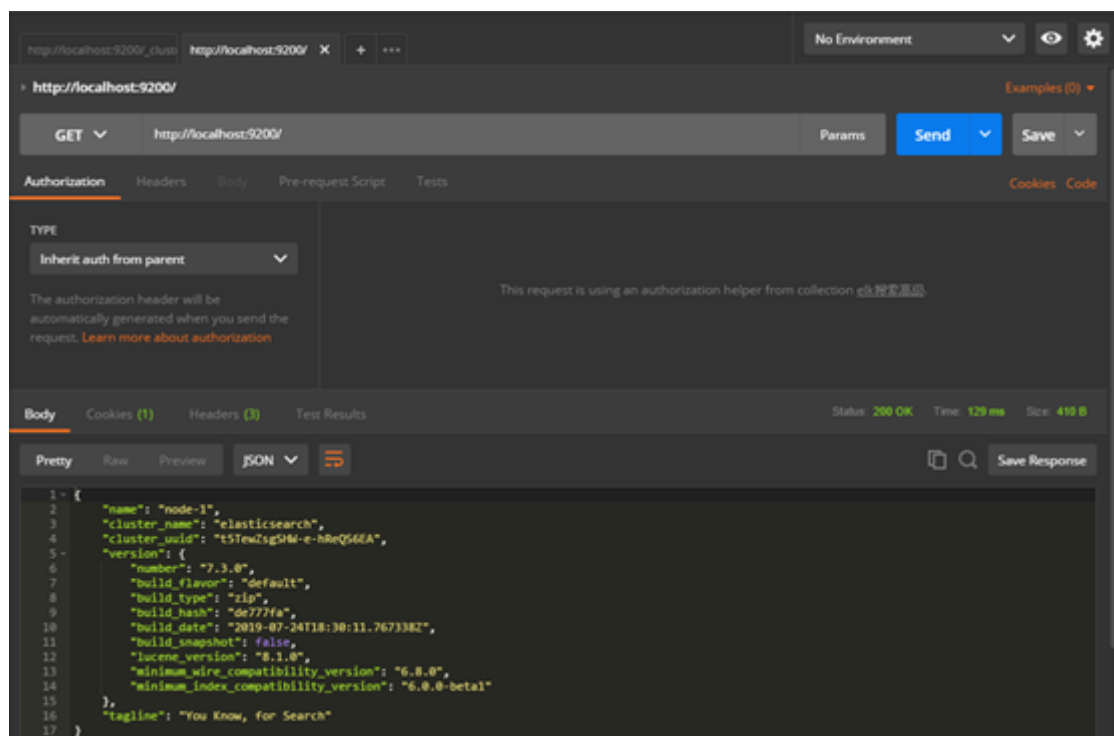
#4.3 Windows安装postman

是什么：postman是一个模拟http请求的工具。能够非常细致地定制化各种http请求。如get\|post\|pu\delete,携带body参数等。

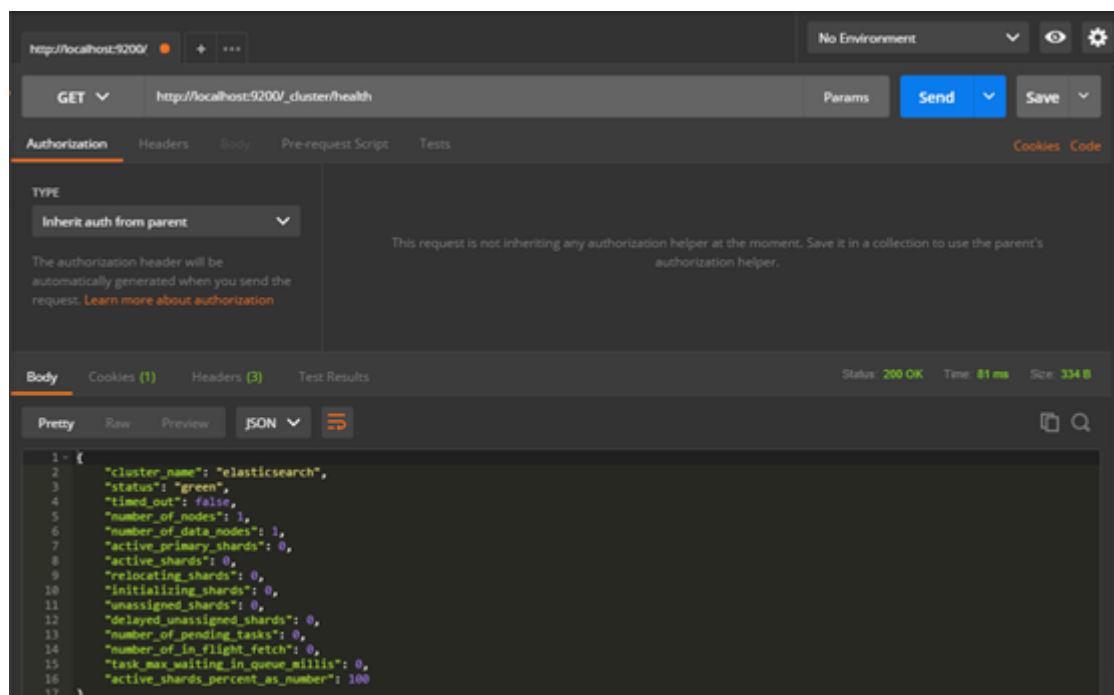
为什么：在没有kibana时，可以使用postman调试。

怎么用：

get <http://localhost:9200/>



测试一下get方式查询集群状态 http://localhost:9200/_cluster/health



#4.4 Windows安装head插件

head插件是ES的一个可视化管理插件，用来监视ES的状态，并通过head客户端和ES服务进行交互，比如创建映射、创建索引等，head的项目地址在 <https://github.com/mobz/elasticsearch-head>。

从ES6.0开始，head插件支持使得node.js运行。

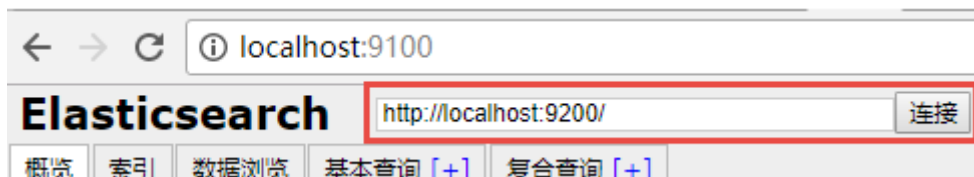
#1安装node.js

#2下载head并运行

```
1 git clone git://github.com/mobz/elasticsearch-head.git
2 cd elasticsearch-head
3 npm install
4 npm run start
```

浏览器打开 <http://localhost:9100/>

#3运行



打开浏览器调试工具发现报错：

Origin null is not allowed by Access-Control-Allow-Origin.

原因是：head插件作为客户端要连接ES服务（localhost:9200），此时存在跨域问题，elasticsearch默认不允许跨域访问。

解决方案：

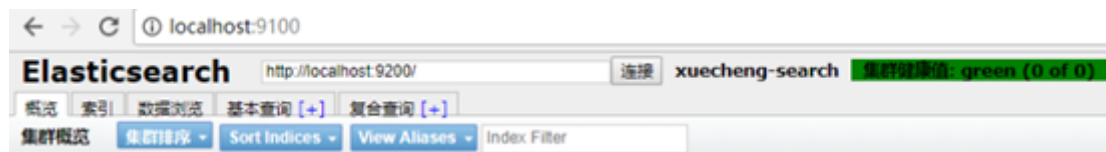
设置elasticsearch允许跨域访问。

在config/elasticsearch.yml 后面增加以下参数：

```
1 #开启cors跨域访问支持，默认为false
2 http.cors.enabled: true
3 #跨域访问允许的域名地址，（允许所有域名）以上使用正则
4 http.cors.allow-origin: /.*/
```

注意：将config/elasticsearch.yml另存为utf-8编码格式。

成功连接ES



注意：kibana\postman\head插件选择自己喜欢的一种使用即可。

本教程使用kibana的dev tool，因为地址栏省略了 http://localhost:9200。

#5. es快速入门

#5.1. 文档（document）的数据格式

- (1) 应用系统的数据结构都是面向对象的，具有复杂的数据结构
- (2) 对象存储到数据库，需要将关联的复杂对象属性插到另一张表，查询时再拼接起来。

(3) es面向文档，文档中存储的数据结构，与对象一致。所以一个对象可以直接存成一个文档。

(4) es的document用json数据格式来表达。

例如：班级和学生关系

```
1  public class Student {
2      private String id;
3      private String name;
4
5      private String classInfoId;
6  }
7
8  private class ClassInfo {
9      private String id;
10     private String className;
11     . . . . .
12
13 }
```

数据库中要设计所谓的一对多，多对一的两张表，外键等。查询出来时，还要关联，mybatis写映射文件，很繁琐。

而在es中，一个学生生成文档如下：

```
1  {
2      "id": "1",
3      "name": "张三",
4      "last_name": "zhang",
5      "classInfo": {
6          "id": "1",
7          "className": "三年二班",
8      }
9  }
```

#5.2图书网站商品管理案例：背景介绍

有一个售卖图书的网站，需要为其基于ES构建一个后台系统，提供以下功能：

- (1) 对商品信息进行CRUD（增删改查）操作
- (2) 执行简单的结构化查询
- (3) 可以执行简单的全文检索，以及复杂的phrase（短语）检索
- (4) 对于全文检索的结果，可以进行高亮显示
- (5) 对数据进行简单的聚合分析

#5.3. 简单的集群管理

5.3.1 快速检查集群的健康状况

es提供了一套api, 叫做cat api, 可以查看es中各种各样的数据

GET /_cat/health?v

1	epoch	timestamp	cluster	status	node.total	node.data	shards	pri	relo
	init	unassign	pending_tasks	max_task_wait_time	active_shards_percent				
2	1568635460	12:04:20	elasticsearch	green	1	1	4	4	0
	0	0	0	-		100.0%			

如何快速了解集群的健康状况? green、yellow、red?

green: 每个索引的primary shard和replica shard都是active状态的

yellow: 每个索引的primary shard都是active状态的, 但是部分replica shard不是active状态, 处于不可用的状态

red: 不是所有索引的primary shard都是active状态的, 部分索引有数据丢失了

5.3.2 快速查看集群中有哪些索引

GET /_cat/indices?v

1	health	status	index	uuid	pri	rep
	docs.count	docs.deleted	store.size	pri.store.size		
2	green	open	.kibana_task_manager	JBMgpuc0SzenstLcjA_G4A	1	0
	2	0	45.5kb	45.5kb		
3	green	open	.monitoring-kibana-7-2019.09.16	LIskf15DTcS70n4Q6t2bTA	1	0
	433	0	218.2kb	218.2kb		
4	green	open	.monitoring-es-7-2019.09.16	RMeUN3tQRjqM8xBgw7Zong	1	0
	3470	1724	1.9mb	1.9mb		
5	green	open	.kibana_1	1cRiyIdATya5xS6qK5pGJw	1	0
	4	0	18.2kb	18.2kb		

5.3.3 简单的索引操作

创建索引: PUT /demo_index?pretty

```
1 {
2   "acknowledged" : true,
3   "shards_acknowledged" : true,
4   "index" : "demo_index"
5 }
```

删除索引: DELETE /demo_index?pretty

5.4 商品的CRUD操作 (document CRUD操作)

5.4.1 新建图书索引

首先建立图书索引 book

语法: put /index

PUT /book



#5.4.2 新增图书:新增文档

语法: PUT /index/type/id

```
1 PUT /book/_doc/1
2
3 {
4   "name": "Bootstrap开发",
5   "description": "Bootstrap是由Twitter推出的一个前台页面开发css框架，是一个非常流行的开发框架，此框架集成了多种页面效果。此开发框架包含了大量的CSS、JS程序代码，可以帮助开发者（尤其是不擅长css页面开发的程序人员）轻松的实现一个css，不受浏览器限制的精美界面css效果。",
6   "studymodel": "201002",
7   "price":38.6,
8   "timestamp":"2019-08-25 19:11:35",
9   "pic":"group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
10  "tags": [ "bootstrap", "dev"]
11 }
```

```
1 PUT /book/_doc/2
2 {
3   "name": "java编程思想",
4   "description": "java语言是世界第一编程语言，在软件开发领域使用人数最多。",
5   "studymodel": "201001",
6   "price":68.6,
7   "timestamp":"2019-08-25 19:11:35",
8   "pic":"group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
9   "tags": [ "java", "dev"]
10 }
```

```
1 PUT /book/_doc/3
2 {
3   "name": "spring开发基础",
4   "description": "spring 在java领域非常流行，java程序员都在用。",
5   "studymodel": "201001",
6   "price":88.6,
7   "timestamp":"2019-08-24 19:11:35",
8   "pic":"group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
9   "tags": [ "spring", "java"]
10 }
```

结果

```
1 {
2   "_index" : "book",
3   "_type" : "_doc",
4   "_id" : "1",
5   "_version" : 1,
6   "result" : "created",
7   "_shards" : {
8     "total" : 2,
9     "successful" : 1,
10    "failed" : 0
11 }
```

```
11     },
12     "_seq_no" : 0,
13     "_primary_term" : 1
14 }
```

#5.4.3 查询图书：检索文档

语法: GET /index/type/id

查看图书:GET /book/_doc/1 就可看到json形式的文档。方便程序解析。

```
1  {
2    "_index" : "book",
3    "_type" : "_doc",
4
5    "_id" : "1",
6
7    "_version" : 4,
8
9    "_seq_no" : 5,
10
11    "_primary_term" : 1,
12
13    "found" : true,
14
15    "_source" : {
16
17      "name" : "Bootstrap开发",
18
19      "description" : "Bootstrap是由Twitter推出的一个前台页面开发css框架，是一个非常流行的开发框架，此框架集成了多种页面效果。此开发框架包含了大量的CSS、JS程序代码，可以帮助开发者（尤其是不擅长css页面开发的程序人员）轻松的实现一个css，不受浏览器限制的精美界面css效果。",
20
21      "studymodel" : "201002",
22
23      "price" : 38.6,
24
25      "timestamp" : "2019-08-25 19:11:35",
26
27      "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
28
29      "tags" : [
30
31        "bootstrap",
32
33        "开发"
34
35      ]
36    }
37  }
38
39 }
```

为方便查看索引中的数据，kibana可以如下操作

Kibana-discover- Create index pattern- Index pattern填book

Create index pattern

Kibana uses index patterns to retrieve data from Elasticsearch indices for things like visualizations.

☐ Include system indices

Step 1 of 2: Define index pattern

Index pattern

book

You can use a * as a wildcard in your index pattern.
You can't use spaces or the characters [, , \ , / , < , > , |].

✓ Success! Your index pattern matches 1 index.

book

Rows per page: 10

> Next step

下一步，再点击discover就可看到数据。

The screenshot shows the Kibana Discover interface. On the left, the 'book' index pattern is selected. Below it, a list of fields is shown, including _source, _id, _index, _score, _type, description, name, price, studymodel, tags, and timestamp. The main area displays a table of documents. The first document has the following fields: name (java编程思想), description (java语言是世界上最流行、在软件开发领域使用人数最多), price (68.6), and timestamp (2019-08-25 19:11:35). A 'JSON' tab is visible, showing the raw document data.

点击json还可以看到原始数据

The screenshot shows the Kibana Discover interface. On the left, the 'book' index pattern is selected. Below it, a list of fields is shown, including _source, _id, _index, _score, _type, description, name, price, studymodel, tags, and timestamp. The main area displays a table of documents. The first document has the following fields: name (java编程思想), description (java语言是世界上最流行、在软件开发领域使用人数最多), price (68.6), and timestamp (2019-08-25 19:11:35). A 'JSON' tab is visible, showing the raw document data.

为方便查看索引中的数据，head可以如下操作

点击数据浏览，点击book索引。

The screenshot shows the Elasticsearch head interface. On the left, the 'book' index pattern is selected. Below it, a list of fields is shown, including _source, _id, _index, _score, _type, description, name, price, studymodel, tags, and timestamp. The main area displays a table of documents. The first document has the following fields: name (java编程思想), description (java语言是世界上最流行、在软件开发领域使用人数最多), price (68.6), and timestamp (2019-08-25 19:11:35). A 'JSON' tab is visible, showing the raw document data.

#5.4.4 修改图书：替换操作

```
1  PUT /book/_doc/1
2  {
3      "name": "Bootstrap开发教程1",
4      "description": "Bootstrap是由Twitter推出的一个前台页面开发css框架，是一个非常流行的开发框架，此框架集成了多种页面效果。此开发框架包含了大量的CSS、JS程序代码，可以帮助开发者（尤其是不擅长css页面开发的程序人员）轻松的实现一个css，不受浏览器限制的精美界面css效果。",
5      "studymodel": "201002",
6      "price":38.6,
7      "timestamp":"2019-08-25 19:11:35",
8      "pic":"group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
9      "tags": [ "bootstrap", "开发" ]
10 }
```

替换操作是整体覆盖，要带上所有信息。

#5.4.5 修改图书：更新文档

语法：POST /{index}/type /{id}/_update

或者POST /{index}/_update/{id}

```
1  POST /book/_update/1/
2  {
3      "doc": {
4          "name": " Bootstrap开发教程高级"
5      }
6  }
```

返回：

```
1  {
2      "_index" : "book",
3      "_type" : "_doc",
4      "_id" : "1",
5      "_version" : 10,
6      "result" : "updated",
7      "_shards" : {
8          "total" : 2,
9          "successful" : 1,
10         "failed" : 0
11     },
12     "_seq_no" : 11,
13     "_primary_term" : 1
14 }
```

#5.4.6 删除图书：删除文档

语法：

```
1  DELETE /book/_doc/1
```

返回：

```
1  {
2      "_index" : "book",
3      "_type" : "_doc",
```

```

4     "_id" : "1",
5     "_version" : 11,
6     "result" : "deleted",
7     "_shards" : {
8         "total" : 2,
9         "successful" : 1,
10        "failed" : 0
11    },
12    "_seq_no" : 12,
13    "_primary_term" : 1
14
15 }

```

#6. 文档document入门

#6.1. 默认自带字段解析

```

1  {
2      "_index" : "book",
3      "_type" : "_doc",
4      "_id" : "1",
5      "_version" : 1,
6      "_seq_no" : 10,
7      "_primary_term" : 1,
8      "found" : true,
9      "_source" : {
10         "name" : "Bootstrap开发教程1",
11         "description" : "Bootstrap是由Twitter推出的一个前台页面开发css框架，是一个非常流行的开发框架，此框架集成了多种页面效果。此开发框架包含了大量的CSS、JS程序代码，可以帮助开发者（尤其是不擅长css页面开发的程序人员）轻松的实现一个css，不受浏览器限制的精美界面css效果。",
12         "studymodel" : "201002",
13         "price" : 38.6,
14         "timestamp" : "2019-08-25 19:11:35",
15         "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
16         "tags" : [
17             "bootstrap",
18             "开发"
19         ]
20     }
21 }

```

#6.1.1 _index

- 含义：此文档属于哪个索引
- 原则：类似数据放在一个索引中。数据库中表的定义规则。如图书信息放在book索引中，员工信息放在employee索引中。各个索引存储和搜索时互不影响。
- 定义规则：英文小写。尽量不要使用特殊字符。order user

#6.1.2 _type

- 含义：类别。book java node
- 注意：以后的es9将彻底删除此字段，所以当前版本在不断弱化type。不需要关注。见到_type都为doc。

#6.1.3 _id

含义：文档的唯一标识。就像表的id主键。结合索引可以标识和定义一个文档。

生成：手动（put /index/_doc/id）、自动

#6.1.4 创建索引时，不同数据放到不同索引中

#6.2. 生成文档id

#6.2.1 手动生成id

场景：数据从其他系统导入时，本身有唯一主键。如数据库中的图书、员工信息等。

用法：put /index/_doc/id

```
1  PUT /test_index/_doc/1
2  {
3    "test_field": "test"
4  }
```

#6.2.2 自动生成id

用法：POST /index/_doc

```
1  POST /test_index/_doc
2  {
3    "test_field": "test1"
4  }
```

返回：

```
1  {
2    "_index" : "test_index",
3    "_type" : "_doc",
4    "_id" : "x29L0m0BPsY0gSJFYZA1",
5    "_version" : 1,
6    "result" : "created",
7    "_shards" : {
8      "total" : 2,
9      "successful" : 1,
10     "failed" : 0
11   },
12   "_seq_no" : 0,
13   "_primary_term" : 1
14 }
```

自动id特点：

长度为20个字符，URL安全，base64编码，GUID，分布式生成不冲突

#6.3. _source 字段

#6.3.1 _source

含义：插入数据时的所有字段和值。在get获取数据时，在_source字段中原样返回。

GET /book/_doc/1

#6.3.2 定制返回字段

就像sql不要select *,而要select name,price from book ...一样。

GET /book/_doc/1?__source_includes=name,price

```
1  {
2    "_index" : "book",
3    "_type" : "_doc",
4    "_id" : "1",
5    "_version" : 1,
6    "_seq_no" : 10,
7    "_primary_term" : 1,
8    "found" : true,
9    "_source" : {
10     "price" : 38.6,
11     "name" : "Bootstrap开发教程1"
12   }
13 }
```

#6.4. 文档的替换与删除

#6.4.1 全量替换

执行两次，返回结果中版本号（_version）在不断上升。此过程为全量替换。

```
1  PUT /test_index/_doc/1
2  {
3    "test_field": "test"
4  }
```

实质：旧文档的内容不会立即删除，只是标记为deleted。适当的时机，集群会将这些文档删除。

#6.4.2 强制创建

为防止覆盖原有数据，我们在新增时，设置为强制创建，不会覆盖原有文档。

语法：PUT /index/ doc/id/ create

```
1  PUT /test_index/_doc/1/_create
2  {
3    "test_field": "test"
4  }
```

返回

```
1  {
2    "error": {
3      "root_cause": [
4        {
5          "type": "version_conflict_engine_exception",
6          "reason": "[2]: version conflict, document already exists (current
version [1])",
7          "index_uuid": "lqzVqxZLQuCnd6LYtZsMkg",
8          "shard": "0",
9          "index": "test_index"
10         }
11      ],
12      "type": "version_conflict_engine_exception",
13      "reason": "[2]: version conflict, document already exists (current version
[1])",
14      "index_uuid": "lqzVqxZLQuCnd6LYtZsMkg",
15      "shard": "0",
16      "index": "test_index"
17    },
18    "status": 409
19  }
```

#6.4.3 删除

DELETE /index/_doc/id

```
1  DELETE  /test_index/_doc/1/
```

实质：旧文档的内容不会立即删除，只是标记为deleted。适当的时机，集群会将这些文档删除。

lazy delete

#6.5. 局部替换 partial update

使用 PUT /index/type/id 为文档全量替换，需要将文档所有数据提交。

partial update局部替换则只修改变动字段。

用法：

```
1  post /index/type/id/_update
2  {
3    "doc": {
4      "field": "value"
5    }
6  }
```

图解内部原理

内部与全量替换是一样的，旧文档标记为删除，新建一个文档。

优点：

- 大大减少网络传输次数和流量，提升性能
- 减少并发冲突发生的概率。

#演示

插入文档

```
1  PUT /test_index/_doc/5
2  {
3    "test_field1": "ydl",
4    "test_field2": "ydlclass"
5  }
```

修改字段1

```
1  POST /test_index/_doc/5/_update
2  {
3    "doc": {
4      "test_field2": " ydlclass 2"
5    }
6  }
```

#6.6. 使用脚本更新

es可以内置脚本执行复杂操作。例如painless脚本。

注意：groovy脚本在es6以后就不支持了。原因是耗内存，不安全远程注入漏洞。

#6.6.1内置脚本

需求1：修改文档6的num字段，+1。

插入数据

```
1  PUT /test_index/_doc/6
2  {
3    "num": 0,
4    "tags": []
5  }
```

执行脚本操作

```
1  POST /test_index/_doc/6/_update
2  {
3    "script" : "ctx._source.num+=1"
4  }
```

查询数据

```
1  GET /test_index/_doc/6
```

返回


```
1  {
2    "_index" : "test_index",
3    "_type" : "_doc",
4    "_id" : "6",
5    "_version" : 2,
6    "_seq_no" : 23,
7    "_primary_term" : 1,
8    "found" : true,
9    "_source" : {
10     "num" : 1,
11     "tags" : [ ]
12   }
13 }
```

需求2: 搜索所有文档, 将num字段乘以2输出

插入数据

```
1  PUT /test_index/_doc/7
2  {
3    "num": 5
4  }
```

查询

```
1  GET /test_index/_search
2  {
3    "script_fields": {
4      "my_doubled_field": {
5        "script": {
6          "lang": "expression",
7          "source": "doc['num'] * multiplier",
8          "params": {
9            "multiplier": 2
10         }
11       }
12     }
13   }
14 }
```

返回

```
1  {
2    "_index" : "test_index",
3    "_type" : "_doc",
4    "_id" : "7",
5    "_score" : 1.0,
6    "fields" : {
7      "my_doubled_field" : [
8        10.0
9      ]
10   }
11 }
```

#6.6.2 外部脚本

Painless是内置支持的。脚本内容可以通过多种途径传给 es，包括 rest 接口，或者放到 config/scripts目录等，默认开启。

注意：脚本性能低下，且容易发生注入，本教程忽略。

官方文档：<https://www.elastic.co/guide/en/elasticsearch/reference/current/modules-scripting-using.html>

#6.7. 图解es的并发问题

如同秒杀，多线程情况下，es同样会出现并发冲突问题。

#6.8. 图解悲观锁与乐观锁机制

为控制并发问题，我们通常采用锁机制。分为悲观锁和乐观锁两种机制。

悲观锁：很悲观，所有情况都上锁。此时只有一个线程可以操作数据。具体例子为数据库中的行级锁、表级锁、读锁、写锁等。

特点：优点是方便，直接加锁，对程序透明。缺点是效率低。

乐观锁：很乐观，对数据本身不加锁。提交数据时，通过一种机制验证是否存在冲突，如es中通过版本号验证。

特点：优点是并发能力高。缺点是操作繁琐，在提交数据时，可能反复重试多次。

#6.9. 图解es内部基于_version乐观锁控制

实验基于_version的版本控制

es对于文档的增删改都是基于版本号。

1新增多次文档：

```
1  PUT /test_index/_doc/3
2  {
3    "test_field": "test"
4  }
```

返回版本号递增

2删除此文档

```
1  DELETE /test_index/_doc/3
```

返回

```
1  DELETE /test_index/_doc/3
2  {
3    "_index" : "test_index",
4    "_type" : "_doc",
5    "_id" : "2",
```

```
6     "_version" : 6,
7     "result" : "deleted",
8     "_shards" : {
9         "total" : 2,
10        "successful" : 1,
11        "failed" : 0
12    },
13    "_seq_no" : 7,
14    "_primary_term" : 1
15 }
```

3再新增

```
1  PUT /test_index/_doc/3
2  {
3      "test_field": "test"
4  }
```

可以看到版本号依然递增，验证延迟删除策略。

如果删除一条数据立马删除的话，所有分片和副本都要立马删除，对es集群压力太大。

#图解es内部并发控制

es内部主从同步时，是多线程异步。乐观锁机制。

#6.10. 演示客户端程序基于_version并发操作流程

java python客户端更新的机制。

#新建文档

```
1  PUT /test_index/_doc/5
2  {
3      "test_field": "ydl"
4  }
```

返回：

```
1  {
2      "_index" : "test_index",
3      "_type" : "_doc",
4      "_id" : "3",
5      "_version" : 1,
6      "result" : "created",
7      "_shards" : {
8          "total" : 2,
9          "successful" : 1,
10         "failed" : 0
11     },
12     "_seq_no" : 8,
13     "_primary_term" : 1
14 }
```

客户端1修改。带版本号1。

首先获取数据的当前版本号

```
1 GET /test_index/_doc/5
```

更新文档

```
1 PUT /test_index/_doc/5?version=1
2 {
3   "test_field": "yd11"
4 }
5 PUT /test_index/_doc/5?if_seq_no=21&if_primary_term=1
6 {
7   "test_field": "yd11"
8 }
```

客户端2并发修改。带版本号1。

```
1 PUT /test_index/_doc/5?version=1
2 {
3   "test_field": "yd12"
4 }
5 PUT /test_index/_doc/5?if_seq_no=21&if_primary_term=1
6 {
7   "test_field": "yd11"
8 }
```

报错。

客户端2重新查询。得到最新版本为2。seq_no=22

```
1 GET /test_index/_doc/4
```

客户端2并发修改。带版本号2。

```
1 PUT /test_index/_doc/4?version=2
2 {
3   "test_field": "yd12"
4 }
5 es7
6 PUT /test_index/_doc/5?if_seq_no=22&if_primary_term=1
7 {
8   "test_field": "yd12"
9 }
```

修改成功。

#6.11. 演示自己手动控制版本号 external version

背景：已有数据是在数据库中，有自己手动维护的版本号的情况下，可以使用external version控制。hbase。

要求：修改时external version要大于当前文档的_version

对比：基于_version时，修改的文档version等于当前文档的版本号。

使用?version=1&version_type=external

#新建文档

```
1  PUT /test_index/_doc/4
2  {
3    "test_field": "yd1"
4  }
```

更新文档:

#客户端1修改文档

```
1  PUT /test_index/_doc/4?version=2&version_type=external
2  {
3    "test_field": "yd11"
4  }
```

#客户端2同时修改

```
1  PUT /test_index/_doc/4?version=2&version_type=external
2  {
3    "test_field": "yd12"
4  }
```

返回:

```
1  {
2    "error": {
3      "root_cause": [
4        {
5          "type": "version_conflict_engine_exception",
6          "reason": "[4]: version conflict, current version [2] is higher or equal
to the one provided [2]",
7          "index_uuid": "-rqYZ2EcSPqL6pu8Gi35jw",
8          "shard": "1",
9          "index": "test_index"
10         }
11       ],
12       "type": "version_conflict_engine_exception",
13       "reason": "[4]: version conflict, current version [2] is higher or equal to
the one provided [2]",
14       "index_uuid": "-rqYZ2EcSPqL6pu8Gi35jw",
15       "shard": "1",
16       "index": "test_index"
17     },
18     "status": 409
19   }
```

#客户端2重新查询数据

```
1  GET /test_index/_doc/4
```

#客户端2重新修改数据

```
1  PUT /test_index/_doc/4?version=3&version_type=external
2  {
3    "test_field": "yd12"
4  }
```

#6.12. 更新时 retry_on_conflict 参数

#retry_on_conflict

指定重试次数

```
1  POST /test_index/_doc/5/_update?retry_on_conflict=3
2  {
3    "doc": {
4      "test_field": "yd11"
5    }
6  }
```

#与 _version 结合使用

```
1  POST /test_index/_doc/5/_update?
    retry_on_conflict=3&version=22&version_type=external
2  {
3    "doc": {
4      "test_field": "yd11"
5    }
6  }
```

#6.13. 批量查询 mget

单条查询 GET /test_index/_doc/1, 如果查询多个id的文档一条一条查询, 网络开销太大。

#mget 批量查询:

```
1  GET /_mget
2  {
3    "docs" : [
4      {
5        "_index" : "test_index",
6        "_type" : "_doc",
7        "_id" : 1
8      },
9      {
10       "_index" : "test_index",
11       "_type" : "_doc",
12       "_id" : 7
13     }
14   ]
15 }
```

返回:


```

1  {
2    "docs" : [
3      {
4        "_index" : "test_index",
5        "_type" : "_doc",
6        "_id" : "2",
7        "_version" : 6,
8        "_seq_no" : 12,
9        "_primary_term" : 1,
10       "found" : true,
11       "_source" : {
12         "test_field" : "test12333123321321"
13       }
14     },
15     {
16       "_index" : "test_index",
17       "_type" : "_doc",
18       "_id" : "3",
19       "_version" : 6,
20       "_seq_no" : 18,
21       "_primary_term" : 1,
22       "found" : true,
23       "_source" : {
24         "test_field" : "test3213"
25       }
26     }
27   ]
28 }

```

提示去掉type

```

1  GET /_mget
2  {
3    "docs" : [
4      {
5        "_index" : "test_index",
6        "_id" : 2
7      },
8      {
9        "_index" : "test_index",
10       "_id" : 3
11     }
12   ]
13 }

```

同一索引下批量查询：

```
1  GET /test_index/_mget
2  {
3      "docs" : [
4          {
5              "_id" : 2
6          },
7          {
8              "_id" : 3
9          }
10     ]
11 }
```

#第三种写法：搜索写法

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

#6.14. 批量增删改 bulk

Bulk 操作解释将文档的增删改查一些列操作，通过一次请求全都做完。减少网络传输次数。

语法：

```
1  POST /_bulk
2  {"action": {"metadata"}}
3  {"data"}
```

如下操作，删除5，新增14，修改2。

```
1  POST /_bulk
2  { "delete": { "_index": "test_index", "_id": "5" }}
3  { "create": { "_index": "test_index", "_id": "14" }}
4  { "test_field": "test14" }
5  { "update": { "_index": "test_index", "_id": "2" } }
6  { "doc" : {"test_field" : "bulk test" } }
```

总结：

1功能：

- delete：删除一个文档，只要1个json串就可以了
- create：相当于强制创建 PUT /index/type/id/_create
- index：普通的put操作，可以是创建文档，也可以是全量替换文档
- update：执行的是局部更新partial update操作

2格式：每个json不能换行。相邻json必须换行。

3隔离：每个操作互不影响。操作失败的行会返回其失败信息。

4实际用法：bulk请求一次不要太大，否则一下积压到内存中，性能会下降。所以，一次请求几千个操作、大小在几M正好。

#6.15. 文档概念学习总结

章节回顾

1文档的增删改查

2文档字段解析

3内部锁机制

4批量查询修改

es是什么

一个分布式的文档数据存储系统distributed document store。es看做一个分布式nosql数据库。如redis\mongoDB\hbase。

文档数据：es可以存储和操作json文档类型的数据，而且这也是es的核心数据结构。存储系统：es可以对json文档类型的数据进行存储，查询，创建，更新，删除，等等操作。

应用场景

- 大数据。es的分布式特点，水平扩容承载大数据。
- 数据结构灵活。列随时变化。使用关系型数据库将会建立大量的关联表，增加系统复杂度。
- 数据操作简单。就是查询，不涉及事务。

举例

电商页面、传统论坛页面等。面向的对象比较复杂，但是作为终端，没有太复杂的功能（事务），只涉及简单的增删改查crud。

这个时候选用ES这种NoSQL型的数据存储，比传统的复杂的事务强大的关系型数据库，更加合适一些。无论是性能，还是吞吐量，可能都会更好。

#7. Java api 实现文档管理

7.1 es技术特点

1es技术比较特殊，不像其他分布式、大数据课程，hadoop、spark、hbase。es代码层面很好写，难的是概念的理解。

2es最重要的是他的rest api。跨语言的。在真实生产中，探查数据、分析数据，使用rest更方便。

3本课程将会大量讲解内部原理及rest api。java代码会在重要的api后学习。

7.2 java 客户端简单获取数据

java api 文档 <https://www.elastic.co/guide/en/elasticsearch/client/java-rest/7.3/java-rest-overview.html>

low：偏向底层。

high：高级封装。足够。

1导包

```

2      <groupId>org.elasticsearch.client</groupId>
3      <artifactId>elasticsearch-rest-high-level-client</artifactId>
4      <version>7.3.0</version>
5      <exclusions>
6          <exclusion>
7              <groupId>org.elasticsearch</groupId>
8              <artifactId>elasticsearch</artifactId>
9          </exclusion>
10     </exclusions>
11 </dependency>
12 <dependency>
13     <groupId>org.elasticsearch</groupId>
14     <artifactId>elasticsearch</artifactId>
15     <version>7.3.0</version>
16 </dependency>

```

2代码

步骤

1 获取连接客户端

2构建请求

3执行

4获取结果

```

1      //获取连接客户端
2      RestHighLevelClient client = new RestHighLevelClient(
3          RestClient.builder(
4              new HttpHost("localhost", 9200, "http")));
5      //构建请求
6      GetRequest getRequest = new GetRequest("book", "1");
7      // 执行
8      GetResponse getResponse = client.get(getRequest, RequestOptions.DEFAULT);
9      // 获取结果
10     if (getResponse.exists()) {
11         long version = getResponse.getVersion();
12         String sourceAsString = getResponse.getSourceAsString();//检索文档(String形式)
13         System.out.println(sourceAsString);
14     }

```

7.3 结合spring-boot-test测试文档查询

0为什么使用spring boot test

- 当今趋势
- 方便开发
- 创建连接交由spring容器，避免每次请求的网络开销。

1导包

```

1      <dependency>
2          <groupId>org.springframework.boot</groupId>
3          <artifactId>spring-boot-starter</artifactId>
4          <version>2.0.6.RELEASE</version>
5      </dependency>
6      <dependency>
7          <groupId>org.springframework.boot</groupId>
8          <artifactId>spring-boot-starter-test</artifactId>
9          <scope>test</scope>
10         <version>2.0.6.RELEASE</version>
11     </dependency>

```

2配置 application.yml

```

1  spring:
2      application:
3          name: service-search
4  heima:
5      elasticsearch:
6          hostlist: 127.0.0.1:9200 #多个结点中间用逗号分隔

```

3代码

主类

配置类

测试类

```

1  @SpringBootTest
2  @RunWith(SpringRunner.class)
3  // 查询文档
4      @Test
5      public void testGet() throws IOException {
6          // 构建请求
7          GetRequest getRequest = new GetRequest("test_post", "1");
8
9          // ===== 可选参数 start =====
10         // 为特定字段配置 _source_include
11         //      String[] includes = new String[]{"user", "message"};
12         //      String[] excludes = Strings.EMPTY_ARRAY;
13         //      FetchSourceContext fetchSourceContext = new FetchSourceContext(true,
14         //      includes, excludes);
15         //      getRequest.fetchSourceContext(fetchSourceContext);
16
17         // 为特定字段配置 _source_excludes
18         //      String[] includes1 = new String[]{"user", "message"};
19         //      String[] excludes1 = Strings.EMPTY_ARRAY;
20         //      FetchSourceContext fetchSourceContext1 = new FetchSourceContext(true,
21         //      includes1, excludes1);
22         //      getRequest.fetchSourceContext(fetchSourceContext1);
23
24         // 设置路由
25         //      getRequest.routing("routing");
26
27         // ===== 可选参数 end =====
28
29         // 查询 同步查询

```

```

29     GetResponse getResponse = client.get(getRequest, RequestOptions.DEFAULT);
30
31     //异步查询
32     //     ActionListener<GetResponse> listener = new ActionListener<GetResponse>
33     //     {
34     //         //查询成功时的立马执行的方法
35     //         @Override
36     //         public void onResponse(GetResponse getResponse) {
37     //             long version = getResponse.getVersion();
38     //             String sourceAsString = getResponse.getSourceAsString();//检索文
39     //             档(String形式)
40     //             System.out.println(sourceAsString);
41     //         }
42     //         //查询失败时的立马执行的方法
43     //         @Override
44     //         public void onFailure(Exception e) {
45     //             e.printStackTrace();
46     //         }
47     //     };
48     //     //执行异步请求
49     //     client.GetAsync(getRequest, RequestOptions.DEFAULT, listener);
50     //     try {
51     //         Thread.sleep(5000);
52     //     } catch (InterruptedException e) {
53     //         e.printStackTrace();
54     //     }
55
56     // 获取结果
57     if (getResponse.exists()) {
58         long version = getResponse.getVersion();
59
60         String sourceAsString = getResponse.getSourceAsString();//检索文档
61         (String形式)
62         System.out.println(sourceAsString);
63         byte[] sourceAsBytes = getResponse.getSourceAsBytes();//以字节接受
64         Map<String, Object> sourceAsMap = getResponse.getSourceAsMap();
65         System.out.println(sourceAsMap);
66
67     }else {
68
69     }

```

#7.4 结合spring-boot-test测试文档新增

rest api

```

1  PUT test_post/_doc/2
2  {
3      "user": "tomas",
4      "postDate": "2019-07-18",
5      "message": "trying out es1"
6  }

```

代码:

```
1  @Test
2      public void testAdd() throws IOException {
3      //      1构建请求
4          IndexRequest request=new IndexRequest("test_posts");
5          request.id("3");
6      //      =====构建文档=====
7      //      构建方法1
8          String jsonString="{\n" +
9              "    \"user\": \"tomas J\", \n" +
10             "    \"postDate\": \"2019-07-18\", \n" +
11             "    \"message\": \"trying out es3\" \n" +
12             "  }";
13          request.source(jsonString, XContentType.JSON);
14
15      //      构建方法2
16      //      Map<String,Object> jsonMap=new HashMap<>();
17      //      jsonMap.put("user", "tomas");
18      //      jsonMap.put("postDate", "2019-07-18");
19      //      jsonMap.put("message", "trying out es2");
20      //      request.source(jsonMap);
21
22      //      构建方法3
23      //      XContentBuilder builder= XContentFactory.jsonBuilder();
24      //      builder.startObject();
25      //      {
26      //          builder.field("user", "tomas");
27      //          builder.timeField("postDate", new Date());
28      //          builder.field("message", "trying out es2");
29      //      }
30      //      builder.endObject();
31      //      request.source(builder);
32      //      构建方法4
33      //      request.source("user", "tomas",
34      //          "postDate", new Date(),
35      //          "message", "trying out es2");
36      //
37      //      =====可选参数=====
38      //      设置超时时间
39      request.timeout(TimeValue.timeValueSeconds(1));
40      request.timeout("1s");
41
42      //      自己维护版本号
43      //      request.version(2);
44      //      request.versionType(VersionType.EXTERNAL);
45
46
47
48      //      2执行
49      //      同步
50      IndexResponse indexResponse = client.index(request,
RequestOptions.DEFAULT);
51      //      异步
52      //      ActionListener<IndexResponse> listener=new
ActionListener<IndexResponse>() {
53      //          @Override
54      //          public void onResponse(IndexResponse indexResponse) {
```

```

55 //
56 //     }
57 //
58 //     @Override
59 //     public void onFailure(Exception e) {
60 //
61 //     }
62 // };
63 // client.indexAsync(request, RequestOptions.DEFAULT, listener );
64 // try {
65 //     Thread.sleep(5000);
66 // } catch (InterruptedException e) {
67 //     e.printStackTrace();
68 // }
69
70
71 //     3获取结果
72 String index = indexResponse.getIndex();
73 String id = indexResponse.getId();
74 //获取插入的类型
75 if(indexResponse.getResult()== DocWriteResponse.Result.CREATED){
76     DocWriteResponse.Result result=indexResponse.getResult();
77     System.out.println("CREATED:"+result);
78 }else if(indexResponse.getResult()== DocWriteResponse.Result.UPDATED){
79     DocWriteResponse.Result result=indexResponse.getResult();
80     System.out.println("UPDATED:"+result);
81 }
82
83 ReplicationResponse.ShardInfo shardInfo = indexResponse.getShardInfo();
84 if(shardInfo.getTotal()!=shardInfo.getSuccessful()){
85     System.out.println("处理成功的分片数少于总分片! ");
86 }
87 if(shardInfo.getFailed(>0){
88     for (ReplicationResponse.ShardInfo.Failure
89 failure:shardInfo.getFailures()) {
90         String reason = failure.reason();//处理潜在的失败原因
91         System.out.println(reason);
92     }
93 }

```

7.5结合spring-boot-test测试文档修改

rest api

```

1 post /test_posts/_doc/3/_update
2 {
3     "doc": {
4         "user": "tomas J"
5     }
6 }

```

代码:

```

1 @Test
2 public void testUpdate() throws IOException {

```



```

3      //      1构建请求
4      UpdateRequest request = new UpdateRequest("test_posts", "3");
5      Map<String, Object> jsonMap = new HashMap<>();
6      jsonMap.put("user", "tomas JJ");
7      request.doc(jsonMap);
8      //=====可选参数=====
9      request.timeout("1s");//超时时间
10
11      //重试次数
12      request.retryOnConflict(3);
13
14      //设置在继续更新之前，必须激活的分片数
15      //      request.waitForActiveShards(2);
16      //所有分片都是active状态，才更新
17      //      request.waitForActiveShards(ActiveShardCount.ALL);
18
19      //      2执行
20      //      同步
21      UpdateResponse updateResponse = client.update(request,
RequestOptions.DEFAULT);
22      //      异步
23
24      //      3获取数据
25      updateResponse.getId();
26      updateResponse.getIndex();
27
28      //判断结果
29      if (updateResponse.getResult() == DocWriteResponse.Result.CREATED) {
30          DocWriteResponse.Result result = updateResponse.getResult();
31          System.out.println("CREATED:" + result);
32      } else if (updateResponse.getResult() == DocWriteResponse.Result.UPDATED)
{
33          DocWriteResponse.Result result = updateResponse.getResult();
34          System.out.println("UPDATED:" + result);
35      }else if(updateResponse.getResult() == DocWriteResponse.Result.DELETED){
36          DocWriteResponse.Result result = updateResponse.getResult();
37          System.out.println("DELETED:" + result);
38      }else if (updateResponse.getResult() == DocWriteResponse.Result.NOOP){
39          //没有操作
40          DocWriteResponse.Result result = updateResponse.getResult();
41          System.out.println("NOOP:" + result);
42      }
43  }

```

7.6结合spring-boot-test测试文档删除

rest api

```
1 DELETE /test_posts/_doc/3
```

代码

```

1      @Test
2      public void testDelete() throws IOException {
3          //      1构建请求
4          DeleteRequest request =new DeleteRequest("test_posts", "3");

```

```

5         //可选参数
6
7
8         //         2执行
9         DeleteResponse deleteResponse = client.delete(request,
RequestOptions.DEFAULT);
10
11
12         //         3获取数据
13         deleteResponse.getId();
14         deleteResponse.getIndex();
15
16         DocWriteResponse.Result result = deleteResponse.getResult();
17         System.out.println(result);
18     }

```

#7.7结合spring-boot-test测试文档bulk

rest api

```

1  POST /_bulk
2  {"action": {"metadata"}}
3  {"data"}

```

代码

```

1  @Test
2      public void testBulk() throws IOException {
3      //         1创建请求
4          BulkRequest request = new BulkRequest();
5          //         request.add(new IndexRequest("post").id("1").source(XContentType.JSON,
"field", "1"));
6          //         request.add(new IndexRequest("post").id("2").source(XContentType.JSON,
"field", "2"));
7
8          request.add(new UpdateRequest("post", "2").doc(XContentType.JSON, "field",
"3"));
9          request.add(new DeleteRequest("post").id("1"));
10
11         //         2执行
12         BulkResponse bulkResponse = client.bulk(request, RequestOptions.DEFAULT);
13
14         for (BulkItemResponse itemResponse : bulkResponse) {
15             DocWriteResponse itemResponseResponse = itemResponse.getResponse();
16
17             switch (itemResponse.getOpType()) {
18                 case INDEX:
19                 case CREATE:
20                     IndexResponse indexResponse = (IndexResponse)
itemResponseResponse;
21                     indexResponse.getId();
22                     System.out.println(indexResponse.getResult());
23                     break;
24                 case UPDATE:
25                     UpdateResponse updateResponse = (UpdateResponse)
itemResponseResponse;

```

```

26         updateResponse.getIndex();
27         System.out.println(updateResponse.getResult());
28         break;
29         case DELETE:
30             DeleteResponse deleteResponse = (DeleteResponse)
itemResponseResponse;
31             System.out.println(deleteResponse.getResult());
32             break;
33         }
34     }
35 }
36

```

#8. 图解es内部机制

#8.1. 图解es分布式基础

#8.1.1es对复杂分布式机制的透明隐藏特性

- 分布式机制：分布式数据存储及共享。
- 分片机制：数据存储到哪个分片，副本数据写入。
- 集群发现机制：cluster discovery。新启动es实例，自动加入集群。
- shard负载均衡：大量数据写入及查询，es会将数据平均分配。
- shard副本：新增副本数，分片重分配。

#8.1.2Elasticsearch的垂直扩容与水平扩容

垂直扩容：使用更加强大的服务器替代老服务器。但单机存储及运算能力有上线。且成本直线上升。如10t服务器1万。单个10T服务器可能20万。

水平扩容：采购更多服务器，加入集群。大数据。

#8.1.3增减或减少节点时的数据rebalance

新增或减少es实例时，es集群会将数据重新分配。

#8.1.4master节点

功能：

- 创建删除节点
- 创建删除索引

#8.1.5节点对等的分布式架构

- 节点对等，每个节点都能接收所有的请求
- 自动请求路由
- 响应收集

#8.2. 图解分片shard、副本replica机制

8.2.1 shard&replica机制

- (1) 每个index包含一个或多个shard
 - (2) 每个shard都是一个最小工作单元，承载部分数据，lucene实例，完整的建立索引和处理请求的能力
 - (3) 增减节点时，shard会自动在nodes中负载均衡
 - (4) primary shard和replica shard，每个document肯定只存在于某一个primary shard以及其对应的replica shard中，不可能存在于多个primary shard
 - (5) replica shard是primary shard的副本，负责容错，以及承担读请求负载
 - (6) primary shard的数量在创建索引的时候就固定了，replica shard的数量可以随时修改
 - (7) primary shard的默认数量是1，replica默认是1，默认共有2个shard，1个primary shard，1个replica shard
- 注意：es7以前primary shard的默认数量是5，replica默认是1，默认有10个shard，5个primary shard，5个replica shard
- (8) primary shard不能和自己的replica shard放在同一个节点上（否则节点宕机，primary shard和副本都丢失，起不到容错的作用），但是可以和其他primary shard的replica shard放在同一个节点上

8.3图解单node环境下创建index是什么样子的

- (1) 单node环境下，创建一个index，有3个primary shard，3个replica shard
- (2) 集群status是yellow
- (3) 这个时候，只会将3个primary shard分配到仅有的一个node上去，另外3个replica shard是无法分配的
- (4) 集群可以正常工作，但是一旦出现节点宕机，数据全部丢失，而且集群不可用，无法承接任何请求

```
1  PUT /test_index1
2  {
3    "settings" : {
4      "number_of_shards" : 3,
5      "number_of_replicas" : 1
6    }
7  }
```

8.4图解2个node环境下replica shard是如何分配的

- (1) replica shard分配：3个primary shard，3个replica shard，1 node
- (2) primary ---> replica同步
- (3) 读请求：primary/replica

8.5图解横向扩容

- 分片自动负载均衡，分片向空闲机器转移。
- 每个节点存储更少分片，系统资源给与每个分片的资源更多，整体集群性能提高。
- 扩容极限：节点数大于整体分片数，则必有空闲机器。
- 超出扩容极限时，可以增加副本数，如设置副本数为2，总共3*3=9个分片。9台机器同时运行，存储和搜索性能更强。容错性更好。

- 容错性：只要一个索引的所有主分片在，集群就可以运行。

#8.6 图解es容错机制 master选举， replica容错， 数据恢复

以3分片，2副本数，3节点为例介绍。

- master node宕机，自动master选举，集群为red
- replica容错：新master将replica提升为primary shard，yellow
- 重启宕机node，master copy replica到该node，使用原有的shard并同步宕机后的修改，green

#9. 图解文档存储机制

#9.1. 数据路由

#9.1.1文档存储如何路由到相应分片

一个文档，最终会落在主分片的一个分片上，到底应该在哪一个分片？这就是数据路由。

#9.1.2路由算法

```
1 shard = hash(routing) % number_of_primary_shards
```

哈希值对主分片数取模。

举例：

对一个文档进行crud时，都会带一个路由值 routing number。默认为文档_id（可能是手动指定，也可能是自动生成）。

存储1号文档，经过哈希计算，哈希值为2,此索引有3个主分片，那么计算 $2\%3=2$ ，就算出此文档在P2分片上。

决定一个document在哪个shard上，最重要的一个值就是routing值，默认是_id，也可以手动指定，相同的routing值，每次过来，从hash函数中，产出的hash值一定是相同的

无论hash值是几，无论是什么数字，对number_of_primary_shards求余数，结果一定是在 $0\sim\text{number_of_primary_shards}-1$ 之间这个范围内的。0,1,2。

#9.1.3手动指定 routing number

```
1 PUT /test_index/_doc/15?routing=num
2 {
3   "num": 0,
4   "tags": []
5 }
```

场景：在程序中，架构师可以手动指定已有数据的一个属性为路由值，好处是可以定制一类文档数据存储到一个分片中。缺点是设计不好，会造成数据倾斜。

所以，不同文档尽量放到不同的索引中。剩下的事情交给es集群自己处理。

#9.1.4主分片数量不可变

涉及到以往数据的查询搜索，所以一旦建立索引，主分片数不可变。

#9.2. 图解文档的增删改内部机制

增删改可以看做update,都是对数据的改动。一个改动请求发送到es集群，经历以下四个步骤：

- (1) 客户端选择一个node发送请求过去，这个node就是coordinating node（协调节点）
- (2) coordinating node，对document进行路由，将请求转发给对应的node（有primary shard）
- (3) 实际的node上的primary shard处理请求，然后将数据同步到replica node。
- (4) coordinating node，如果发现primary node和所有replica node都搞定之后，就返回响应结果给客户端。

#9.3. 图解文档的查询内部机制

- 1、客户端发送请求到任意一个node，成为coordinate node
- 2、coordinate node对document进行路由，将请求转发到对应的node，此时会使用round-robin随机轮询算法，在primary shard以及其所有replica中随机选择一个，让读请求负载均衡
- 3、接收请求的node返回document给coordinate node
- 4、coordinate node返回document给客户端
- 5、特殊情况：document如果还在建立索引过程中，可能只有primary shard有，任何一个replica shard都没有，此时可能会导致无法读取到document，但是document完成索引建立之后，primary shard和replica shard就都有了。

#9.4. bulk api奇特的json格式

```
1  POST /_bulk
2  {"action": {"meta"}}\n
3  {"data"}\n
4  {"action": {"meta"}}\n
5  {"data"}\n
6
7  [
8    {
9      "action": {
10        "method": "create"
11      },
12      "data": {
13        "id": 1,
14        "field1": "java",
15        "field1": "spring",
16      }
17    },
18    {
19      "action": {
```

```

20         "method": "create"
21     },
22     "data": {
23         "id": 2,
24         "field1": "java",
25         "field1": "spring",
26     }
27 }
28 ]

```

1、bulk中的每个操作都可能要转发到不同的node的shard去执行

2、如果采用比较良好的json数组格式

允许任意的换行，整个可读性非常棒，读起来很爽，es拿到那种标准格式的json串以后，要按照下述流程去进行处理

(1) 将json数组解析为JSONArray对象，这个时候，整个数据，就会在内存中出现一份一模一样的拷贝，一份数据是json文本，一份数据是JSONArray对象

(2) 解析json数组里的每个json，对每个请求中的document进行路由

(3) 为路由到同一个shard上的多个请求，创建一个请求数组。100请求中有10个是到P1.

(4) 将这个请求数组序列化

(5) 将序列化后的请求数组发送到对应的节点上去

3、耗费更多内存，更多的jvm gc开销

我们之前提到过bulk size最佳大小的那个问题，一般建议说在几千条那样，然后大小在10MB左右，所以说，可怕的事情来了。假设说现在100个bulk请求发送到了一个节点上去，然后每个请求是10MB，100个请求，就是1000MB = 1GB，然后每个请求的json都copy一份为jsonarray对象，此时内存中的占用就会翻倍，就会占用2GB的内存，甚至还不止。因为弄成jsonarray之后，还可能会多搞一些其他的数据结构，2GB+的内存占用。

占用更多的内存可能会积压其他请求的内存使用量，比如说最重要的搜索请求，分析请求，等等，此时就可能会导致其他请求的性能急速下降。

另外的话，占用内存更多，就会导致java虚拟机的垃圾回收次数更多，跟频繁，每次要回收的垃圾对象更多，耗费的时间更多，导致es的java虚拟机停止工作线程的时间更多。

4、现在的奇特格式

```

1  POST /_bulk
2  { "delete": { "_index": "test_index", "_id": "5" } } \n
3  { "create": { "_index": "test_index", "_id": "14" } } \n
4  { "test_field": "test14" } \n
5  { "update": { "_index": "test_index", "_id": "2" } } \n
6  { "doc" : { "test_field" : "bulk test" } } \n

```

(1) 不用将其转换为json对象，不会出现内存中的相同数据的拷贝，直接按照换行符切割json

(2) 对每两个一组的json，读取meta，进行document路由

(3) 直接将对应的json发送到node上去

5、最大的优势在于，不需要将json数组解析为一个JSONArray对象，形成一份大数据的拷贝，浪费内存空间，尽可能地保证性能。

#10. Mapping映射入门

#10.1. 什么是mapping映射

概念：自动或手动为index中的_doc建立的一种数据结构和相关配置，简称为mapping映射。

插入几条数据，让es自动为我们建立一个索引

```
1  PUT /website/_doc/1
2  {
3    "post_date": "2019-01-01",
4    "title": "my first article",
5    "content": "this is my first article in this website",
6    "author_id": 11400
7  }
8
9  PUT /website/_doc/2
10 {
11   "post_date": "2019-01-02",
12   "title": "my second article",
13   "content": "this is my second article in this website",
14   "author_id": 11400
15 }
16
17 PUT /website/_doc/3
18 {
19   "post_date": "2019-01-03",
20   "title": "my third article",
21   "content": "this is my third article in this website",
22   "author_id": 11400
23 }
```

对比数据库建表语句

```
1  create table website(
2    post_date date,
3    title varchar(50),
4    content varchar(100),
5    author_id int(11)
6  );
```

动态映射：dynamic mapping，自动为我们建立index，以及对应的mapping，mapping中包含了每个field对应的数据类型，以及如何分词等设置。

重点：我们当然，后面会讲解，也可以手动在创建数据之前，先创建index，以及对应的mapping

```
1  GET /website/_mapping/
2  {
3    "website" : {
4      "mappings" : {
5        "properties" : {
6          "author_id" : {
7            "type" : "long"
8          },
9          "content" : {
```



```

10         "type" : "text",
11         "fields" : {
12             "keyword" : {
13                 "type" : "keyword",
14                 "ignore_above" : 256
15             }
16         },
17     },
18     "post_date" : {
19         "type" : "date"
20     },
21     "title" : {
22         "type" : "text",
23         "fields" : {
24             "keyword" : {
25                 "type" : "keyword",
26                 "ignore_above" : 256
27             }
28         }
29     }
30 }
31 }
32 }
33 }

```

尝试各种搜索

1	GET /website/_search?q=2019	0条结果
2	GET /website/_search?q=2019-01-01	1条结果
3	GET /website/_search?q=post_date:2019-01-01	1条结果
4	GET /website/_search?q=post_date:2019	0 条结果

搜索结果为什么不一致，因为es自动建立mapping的时候，设置了不同的field不同的data type。不同的data type的分词、搜索等行为是不一样的。所以出现了_all field和post_date field的搜索表现完全不一样。

10.2. 精确匹配与全文搜索的对比分析

10.2.1 exact value 精确匹配

2019-01-01, exact value, 搜索的时候，必须输入2019-01-01，才能搜索出来

如果你输入一个01，是搜索不出来的

```
select * from book where name= 'java'
```

10.2.2 full text 全文检索

搜“笔记本电脑”，笔记本电脑词条会不会出现。

```
select * from book where name like '%java%'
```

(1) 缩写 vs. 全称: cn vs. china

(2) 格式转化: like liked likes

(3) 大小写: Tom vs tom

(4) 同义词: like vs love

2019-01-01, 2019 01 01, 搜索2019, 或者01, 都可以搜索出来

china, 搜索cn, 也可以将china搜索出来

likes, 搜索like, 也可以将likes搜索出来

Tom, 搜索tom, 也可以将Tom搜索出来

like, 搜索love, 同义词, 也可以将like搜索出来

就不是说单纯的只是匹配完整的一个值, 而是可以对值进行拆分词语后(分词)进行匹配, 也可以通过缩写、时态、大小写、同义词等进行匹配。深入 NPL,自然语义处理。

10.3. 全文检索下倒排索引核心原理快速揭秘

doc1: I really liked my small dogs, and I think my mom also liked them.

doc2: He never liked any dogs, so I hope that my mom will not expect me to liked him.

分词, 初步的倒排索引的建立

term	doc1	doc2
I	*	*
really	*	
liked	*	*
my	*	*
small	*	
dogs	*	
and	*	
think	*	
mom	*	*
also	*	
them	*	
He		*
never		*
any		*
so		*
hope		*
that		*
will		*
not		*
expect		*
me		*
to		*
him		*

演示了一下倒排索引最简单的建立的一个过程

搜索

mother like little dog, 不可能有任何结果

mother

like

little

dog

这不是我们想要的结果。同义词mom\mother在我们人类看来是一样。想进行标准化操作。

重建倒排索引

normalization正规化，建立倒排索引的时候，会执行一个操作，也就是说对拆分出的各个单词进行相应的处理，以提升后面搜索的时候能够搜索到相关联的文档的概率

时态的转换，单复数的转换，同义词的转换，大小写的转换

mom → mother

liked → like

small → little

dogs → dog

重新建立倒排索引，加入normalization，再次用mother liked little dog搜索，就可以搜索到了

word	doc1	doc2	normalization
I	*	*	
really	*		
like	*	*	liked → like
my	*	*	
little	*		small → little
dog	*		dogs → dog
and	*		
think	*		
mother	*	*	mom → mother
also	*		
them	*		
He		*	
never		*	
any		*	
so		*	
hope		*	
that		*	
will		*	
not		*	
expect		*	
me		*	
to		*	
him		*	

重新搜索

搜索：mother liked little dog,

对搜索条件进行分词 normalization

mother

liked -> like

little

dog

doc1和doc2都会搜索出来

10.4. 分词器 analyzer

10.4.1 什么是分词器 analyzer

作用：切分词语，normalization（提升recall召回率）

给你一段句子，然后将这段句子拆分成一个一个的单个的单词，同时对每个单词进行normalization（形态转换，单复数转换）

recall，召回率：搜索的时候，增加能够搜索到的结果的数量

analyzer 组成部分：

1、character filter：在一段文本进行分词之前，先进行预处理，比如说最常见的就是，过滤html标签（hello --> hello），& --> and（I&you --> I and you）

2、tokenizer：分词，hello you and me --> hello, you, and, me

3、token filter：lowercase，stop word，synonymom，dogs --> dog，liked --> like，Tom --> tom，a/the/an --> 干掉，mother --> mom，small --> little

stop word 停用词：了 的 呢。

一个分词器，很重要，将一段文本进行各种处理，最后处理好的结果才会拿去建立倒排索引。

10.4.2 内置分词器的介绍

例句：Set the shape to semi-transparent by calling set_trans(5)

standard analyzer标准分词器：set, the, shape, to, semi, transparent, by, calling, set_trans, 5（默认的是 standard）

simple analyzer简单分词器：set, the, shape, to, semi, transparent, by, calling, set, trans

whitespace analyzer：Set, the, shape, to, semi-transparent, by, calling, set_trans(5)

language analyzer（特定的语言的分词器，比如说，english，英语分词器）：set, shape, semi, transpar, call, set_tran, 5

官方文档：

<https://www.elastic.co/guide/en/elasticsearch/reference/7.4/analysis-analyzers.html>

Analyzers



Elasticsearch ships with a wide range of built-in analyzers, which can be used in any index without further configuration:

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.

Stop Analyzer

The `stop` analyzer is like the `simple` analyzer, but also supports removal of stop words.

Keyword Analyzer

The `keyword` analyzer is a “noop” analyzer that accepts whatever text it is given and outputs the exact same text as a single term.

Pattern Analyzer

The `pattern` analyzer uses a regular expression to split the text into terms. It supports lower-casing and stop words.

Language Analyzers

Elasticsearch provides many language-specific analyzers like `english` or `french`.

Fingerprint Analyzer

The `fingerprint` analyzer is a specialist analyzer which creates a fingerprint which can be used for duplicate detection.

10.5. query string根据字段分词策略

10.5.1 query string分词

query string必须以和index建立时相同的analyzer进行分词

query string对exact value和full text的区别对待

如: date: exact value 精确匹配

text: full text 全文检索

10.5.2 测试分词器

```
1 GET /_analyze
2 {
3   "analyzer": "standard",
4   "text": "Text to analyze 80"
5 }
```

返回值:

```

1  {
2    "tokens" : [
3      {
4        "token" : "text",
5        "start_offset" : 0,
6        "end_offset" : 4,
7        "type" : "<ALPHANUM>",
8        "position" : 0
9      },
10     {
11       "token" : "to",
12       "start_offset" : 5,
13       "end_offset" : 7,
14       "type" : "<ALPHANUM>",
15       "position" : 1
16     },
17     {
18       "token" : "analyze",
19       "start_offset" : 8,
20       "end_offset" : 15,
21       "type" : "<ALPHANUM>",
22       "position" : 2
23     },
24     {
25       "token" : "80",
26       "start_offset" : 16,
27       "end_offset" : 18,
28       "type" : "<NUM>",
29       "position" : 3
30     }
31   ]
32 }

```

token 实际存储的term 关键字

position 在此词条在原文本中的位置

start_offset/end_offset字符在原始字符串中的位置

#10.6. mapping回顾总结

- (1) 往es里面直接插入数据，es会自动建立索引，同时建立对应的mapping。(dynamic mapping)
- (2) mapping中就自动定义了每个field的数据类型
- (3) 不同的数据类型（比如说text和date），可能有的是exact value，有的是full text
- (4) exact value，在建立倒排索引的时候，分词的时候，是将整个值一起作为一个关键词建立到倒排索引中的；full text，会经历各种各样的处理，分词，normalization（时态转换，同义词转换，大小写转换），才会建立到倒排索引中。
- (5) 同时呢，exact value和full text类型的field就决定了，在一个搜索过来的时候，对exact value field或者是full text field进行搜索的行为也是不一样的，会跟建立倒排索引的行为保持一致；比如说exact value搜索的时候，就是直接按照整个值进行匹配，full text query string，也会进行分词和normalization再去倒排索引中去搜索

(6) 可以用es的dynamic mapping, 让其自动建立mapping, 包括自动设置数据类型; 也可以提前手动创建index和mapping, 自己对各个field进行设置, 包括数据类型, 包括索引行为, 包括分词器, 等。

#10.7. mapping的核心数据类型以及dynamic mapping

#10.7.1 核心的数据类型

string :text and keyword

byte, short, integer, long,float, double

boolean

date

详见: <https://www.elastic.co/guide/en/elasticsearch/reference/7.3/mapping-types.html>

下图是ES7.3核心的字段类型如下:

Field datatypes

Elasticsearch supports a number of different datatypes for the fields in a document:

Core datatypes

string

text and keyword

Numeric

long, integer, short, byte, double, float, half_float, scaled_float

Date

date

Date nanoseconds

date_nanos

Boolean

boolean

Binary

binary

Range

integer_range, float_range, long_range, double_range, date_range

Complex datatypes

Object

object for single JSON objects

Nested

nested for arrays of JSON objects

Geo datatypes

Geo-point

geo_point for lat/lon points

10.7.2 dynamic mapping 推测规则

true or false --> boolean

123 --> long

123.45 --> double

2019-01-01 --> date

"hello world" --> text/keyword

10.7.3 查看mapping

GET /index/_mapping/

10.8 手动管理mapping

10.8.1 查询所有索引的映射

GET /_mapping

10.8.2 创建映射！！

创建索引后，应该立即手动创建映射

```
1  PUT book/_mapping
2  {
3      "properties": {
4          "name": {
5              "type": "text"
6          },
7          "description": {
8              "type": "text",
9              "analyzer": "english",
10             "search_analyzer": "english"
11         },
12         "pic": {
13             "type": "text",
14             "index": false
15         },
16         "studymodel": {
17             "type": "text"
18         }
19     }
20 }
```

Text 文本类型

1) analyzer

通过analyzer属性指定分词器。

上边指定了analyzer是指在索引和搜索都使用english，如果单独想定义搜索时使用的分词器则可以通过search_analyzer属性。

2) index

index属性指定是否索引。

默认为index=true，即要进行索引，只有进行索引才可以从索引库搜索到。

但是也有一些内容不需要索引，比如：商品图片地址只被用来展示图片，不进行搜索图片，此时可以将index设置为false。

删除索引，重新创建映射，将pic的index设置为false，尝试根据pic去搜索，结果搜索不到数据。

3) store

是否在source之外存储，每个文档索引后会在ES中保存一份原始文档，存放在"source"中，一般情况下不需要设置store为true，因为在source中已经有一份原始文档了。

测试

```
1  PUT book/_mapping
2  {
3      "properties": {
4          "name": {
5              "type": "text"
6          },
7          "description": {
8              "type": "text",
9              "analyzer": "english",
10             "search_analyzer": "english"
11          },
12          "pic": {
13              "type": "text",
14              "index": false
15          },
16          "studymodel": {
17              "type": "text"
18          }
19      }
20  }
```

插入文档：

```
1  PUT /book/_doc/1
2  {
3      "name": "Bootstrap开发框架",
4      "description": "Bootstrap是由Twitter推出的一个前台页面开发框架，在行业之中使用较为广泛。此开发框架包含了大量的CSS、JS程序代码，可以帮助开发者（尤其是不擅长页面开发的程序人员）轻松的实现一个不受浏览器限制的精美界面效果。",
5      "pic": "group1/M00/00/01/wKh1QFqO4MmAOP53AAAcwDwm6SU490.jpg",
6      "studymodel": "201002"
7  }
```

Get /book/_search?q=name:开发

Get /book/_search?q=description:开发

Get /book/_search?q=pic:group1/M00/00/01/wKh1QFqO4MmAOP53AAAcwDwm6SU490.jpg

Get /book/_search?q=studymodel:201002

通过测试发现：name和description都支持全文检索，pic不可作为查询条件。

keyword关键字字段

目前已经取代了"index": false。上边介绍的text文本字段在映射时要设置分词器，keyword字段为关键字字段，通常搜索keyword是按照整体搜索，所以创建keyword字段的索引时是不进行分词的，比如：邮政编码、手机号码、身份证等。keyword字段通常用于过滤、排序、聚合等。

date日期类型

日期类型不用设置分词器。

通常日期类型的字段用于排序。

format

通过format设置日期格式

例子：

下边的设置允许date字段存储年月日时分秒、年月日及毫秒三种格式。

```
{
  "properties": {
    "timestamp": {
      "type": "date",
      "format": "yyyy-MM-dd HH:mm:ss||yyyy-MM-dd"
    }
  }
}
```

插入文档：

Post book/doc/3

```
{
  "name": "spring开发基础",
  "description": "spring 在java领域非常流行，java程序员都在用。",
  "studymodel": "201001",
  "pic": "group1/M00/00/01/wKhIQFqO4MmAOP53AAAcwDwm6SU490.jpg",
  "timestamp": "2018-07-04 18:28:58"
}
```

数值类型

下边是ES支持的数值类型

Numeric datatypes



The following numeric types are supported:

<code>long</code>	A signed 64-bit integer with a minimum value of -2^{63} and a maximum value of $2^{63}-1$.
<code>integer</code>	A signed 32-bit integer with a minimum value of -2^{31} and a maximum value of $2^{31}-1$.
<code>short</code>	A signed 16-bit integer with a minimum value of $-32,768$ and a maximum value of $32,767$.
<code>byte</code>	A signed 8-bit integer with a minimum value of -128 and a maximum value of 127 .
<code>double</code>	A double-precision 64-bit IEEE 754 floating point number, restricted to finite values.
<code>float</code>	A single-precision 32-bit IEEE 754 floating point number, restricted to finite values.
<code>half_float</code>	A half-precision 16-bit IEEE 754 floating point number, restricted to finite values.
<code>scaled_float</code>	A floating point number that is backed by a <code>long</code> , scaled by a fixed <code>double</code> scaling factor.

- 1、尽量选择范围小的类型，提高搜索效率
- 2、对于浮点数尽量用比例因子，比如一个价格字段，单位为元，我们将比例因子设置为100这在ES中会按分存储，映射如下：

```
1  "price": {
2      "type": "scaled_float",
3      "scaling_factor": 100
4  },
```

由于比例因子为100，如果我们输入的价格是23.45则ES中会将23.45乘以100存储在ES中。

如果输入的价格是23.456，ES会将23.456乘以100再取一个接近原始值的数，得出2346。

使用比例因子的好处是整型比浮点型更易压缩，节省磁盘空间。

如果比例因子不适合，则从下表选择范围小的去用：

更新已有映射，并插入文档：

```
1  PUT book/doc/3
2  {
3      "name": "spring开发基础",
4      "description": "spring 在java领域非常流行，java程序员都在用。",
5      "studymodel": "201001",
6      "pic": "group1/M00/00/01/wKh1QFq04MmAOP53AAAcwDwm6SU490.jpg",
7      "timestamp": "2018-07-04 18:28:58",
8      "price": 38.6
9  }
```

10.8.3 修改映射

只能创建index时手动建立mapping，或者新增field mapping，但是不能update field mapping。

因为已有数据按照映射早已分词存储好。如果修改，那这些存量数据怎么办。

新增一个字段mapping

```
1  PUT /book/_mapping/
2  {
3    "properties" : {
4      "new_field" : {
5        "type" : "text",
6        "index": "false"
7      }
8    }
9  }
```

如果修改mapping,会报错

```
1  PUT /book/_mapping/
2  {
3    "properties" : {
4      "studymodel" : {
5        "type" : "keyword"
6      }
7    }
8  }
```

返回：

```
1  {
2    "error": {
3      "root_cause": [
4        {
5          "type": "illegal_argument_exception",
6          "reason": "mapper [studymodel] of different type, current_type [text],
merged_type [keyword]"
7        }
8      ],
9      "type": "illegal_argument_exception",
10     "reason": "mapper [studymodel] of different type, current_type [text],
merged_type [keyword]"
11   },
12   "status": 400
13 }
```

10.8.4 删除映射

通过删除索引来删除映射。

10.9 复杂数据类型

#10.9 .1 multivalue field

```
{ "tags": [ "tag1", "tag2" ] }
```

建立索引时与string是一样的，数据类型不能混

#10.9 .2. empty field

null, [], [null]

#10.9 .3. object field

```
1  PUT /company/_doc/1
2  {
3    "address": {
4      "country": "china",
5      "province": "guangdong",
6      "city": "guangzhou"
7    },
8    "name": "jack",
9    "age": 27,
10   "join_date": "2019-01-01"
11 }
```

address: object类型

查询映射

```
1  GET /company/_mapping
2  {
3    "company" : {
4      "mappings" : {
5        "properties" : {
6          "address" : {
7            "properties" : {
8              "city" : {
9                "type" : "text",
10               "fields" : {
11                 "keyword" : {
12                   "type" : "keyword",
13                   "ignore_above" : 256
14                 }
15               }
16             },
17             "country" : {
18               "type" : "text",
19               "fields" : {
20                 "keyword" : {
21                   "type" : "keyword",
22                   "ignore_above" : 256
23                 }
24               }
25             },
26             "province" : {
27               "type" : "text",
28               "fields" : {
29                 "keyword" : {
30                   "type" : "keyword",
```

```

31         "ignore_above" : 256
32     }
33 }
34 }
35 }
36 },
37 "age" : {
38     "type" : "long"
39 },
40 "join_date" : {
41     "type" : "date"
42 },
43 "name" : {
44     "type" : "text",
45     "fields" : {
46         "keyword" : {
47             "type" : "keyword",
48             "ignore_above" : 256
49         }
50     }
51 }
52 }
53 }
54 }
55 }

```

object

```

1  {
2    "address": {
3      "country": "china",
4      "province": "guangdong",
5      "city": "guangzhou"
6    },
7    "name": "jack",
8    "age": 27,
9    "join_date": "2017-01-01"
10 }

```

底层存储格式

```

1  {
2    "name":          [jack],
3    "age":           [27],
4    "join_date":     [2017-01-01],
5    "address.country": [china],
6    "address.province": [guangdong],
7    "address.city":  [guangzhou]
8  }

```

对象数组:


```
1  {
2    "authors": [
3      { "age": 26, "name": "Jack White"},
4      { "age": 55, "name": "Tom Jones"},
5      { "age": 39, "name": "Kitty Smith"}
6    ]
7  }
```

存储格式:

```
1  {
2    "authors.age":    [26, 55, 39],
3    "authors.name":  [jack, white, tom, jones, kitty, smith]
4  }
```

#11. 索引Index入门

#11.1. 索引管理

#11.1.1. 为什么我们要手动创建索引

直接put数据 PUT index/_doc/1,es会自动生成索引, 并建立动态映射dynamic mapping。

在生产上, 我们需要自己手动建立索引和映射, 为了更好地管理索引。就像数据库的建表语句一样。

#11.1.2. 索引管理

#11.1.2.1 创建索引

创建索引的语法

```
1  PUT /index
2  {
3    "settings": { ... any settings ... },
4    "mappings": {
5      "properties" : {
6        "field1" : { "type" : "text" }
7      }
8    },
9    "aliases": {
10     "default_index": {}
11   }
12 }
```

举例:

```
1  PUT /my_index
2  {
3    "settings": {
4      "number_of_shards": 1,
5      "number_of_replicas": 1
6    },
7    "mappings": {
8      "properties": {
```

```
9      "field1":{
10        "type": "text"
11      },
12      "field2":{
13        "type": "text"
14      }
15    }
16  },
17  "aliases": {
18    "default_index": {}
19  }
20 }
```

索引别名

插入数据

```
1  POST /my_index/_doc/1
2  {
3    "field1": "java",
4    "field2": "js"
5  }
```

查询数据 都可以查到

GET /my_index/_doc/1

GET /default_index/_doc/1

11.1.2.2 查询索引

GET /my_index/_mapping

GET /my_index/_setting

11.1.2.3 修改索引

修改副本数

```
1  PUT /my_index/_settings
2  {
3    "index" : {
4      "number_of_replicas" : 2
5    }
6  }
```

11.1.2.4 删除索引

DELETE /my_index

DELETE /index_one,index_two

DELETE /index_*

DELETE /_all

为了安全起见，防止恶意删除索引，删除时必须指定索引名：

elasticsearch.yml

action.destructive_requires_name: true

11.2. 定制分词器

11.2.1 默认的分词器

standard

分词三个组件, character filter, tokenizer, token filter

standard tokenizer: 以单词边界进行切分

standard token filter: 什么都不做

lowercase token filter: 将所有字母转换为小写

stop token filter (默认被禁用): 移除停用词, 比如a the it等等

11.2.2 修改分词器的设置

启用english停用词token filter

```
1  PUT /my_index
2  {
3    "settings": {
4      "analysis": {
5        "analyzer": {
6          "es_std": {
7            "type": "standard",
8            "stopwords": "_english_"
9          }
10       }
11     }
12   }
13 }
```

测试分词

```
1  GET /my_index/_analyze
2  {
3    "analyzer": "standard",
4    "text": "a dog is in the house"
5  }
6
7  GET /my_index/_analyze
8  {
9    "analyzer": "es_std",
10   "text": "a dog is in the house"
11 }
```

11.2.3 定制化自己的分词器

```
1  PUT /my_index
2  {
3    "settings": {
4      "analysis": {
5        "char_filter": {
6          "&_to_and": {
7            "type": "mapping",
```

```

8         "mappings": ["&=> and"]
9     }
10 },
11     "filter": {
12         "my_stopwords": {
13             "type": "stop",
14             "stopwords": ["the", "a"]
15         }
16     },
17     "analyzer": {
18         "my_analyzer": {
19             "type": "custom",
20             "char_filter": ["html_strip", "&_to_and"],
21             "tokenizer": "standard",
22             "filter": ["lowercase", "my_stopwords"]
23         }
24     }
25 }
26 }
27 }

```

测试

```

1 GET /my_index/_analyze
2 {
3     "analyzer": "my_analyzer",
4     "text": "tom&jerry are a friend in the house, <a>, HAHA!!"
5 }

```

设置字段使用自定义分词器

```

1 PUT /my_index/_mapping/
2 {
3     "properties": {
4         "content": {
5             "type": "text",
6             "analyzer": "my_analyzer"
7         }
8     }
9 }

```

11.3 type底层结构及弃用原因

11.3.1 type是什么

type，是一个index中用来区分类似的数据的，类似的数据，但是可能有不同的fields，而且有不同的属性来控制索引建立、分词器。field的value，在底层的lucene中建立索引的时候，全部是opaque bytes类型，不区分类型的。lucene是没有type的概念的，在document中，实际上将type作为一个document的field来存储，即 type，es通过 type来进行type的过滤和筛选。

11.3.2es中不同type存储机制

一个index中的多个type，实际上是放在一起存储的，因此一个index下，不能有多个type重名，而类型或者其他设置不同的，因为那样是无法处理的

```
1  {
2    "goods": {
3      "mappings": {
4        "electronic_goods": {
5          "properties": {
6            "name": {
7              "type": "string",
8            },
9            "price": {
10             "type": "double"
11           },
12           "service_period": {
13             "type": "string"
14           }
15         }
16       },
17       "fresh_goods": {
18         "properties": {
19           "name": {
20             "type": "string",
21           },
22           "price": {
23             "type": "double"
24           },
25           "eat_period": {
26             "type": "string"
27           }
28         }
29       }
30     }
31   }
32 }
```

```
1  PUT /goods/electronic_goods/1
2  {
3    "name": "小米空调",
4    "price": 1999.0,
5    "service_period": "one year"
6  }
```

```
1  PUT /goods/fresh_goods/1
2  {
3    "name": "澳洲龙虾",
4    "price": 199.0,
5    "eat_period": "one week"
6  }
```

es文档在底层的存储是这样子的

```
1  {
2    "goods": {
3      "mappings": {
```

```

4      "_type": {
5        "type": "string",
6        "index": "false"
7      },
8      "name": {
9        "type": "string"
10     }
11     "price": {
12       "type": "double"
13     }
14     "service_period": {
15       "type": "string"
16     },
17     "eat_period": {
18       "type": "string"
19     }
20   }
21 }
22 }

```

底层数据存储格式

```

1  {
2    "_type": "electronic_goods",
3    "name": "小米空调",
4    "price": 1999.0,
5    "service_period": "one year",
6    "eat_period": ""
7  }

```

```

1  {
2    "_type": "fresh_goods",
3    "name": "澳洲龙虾",
4    "price": 199.0,
5    "service_period": "",
6    "eat_period": "one week"
7  }
8

```

11.3.3 type弃用

同一索引下，不同type的数据存储其他type的field 大量空值，造成资源浪费。

所以，不同类型数据，要放到不同的索引中。

es9中，将会彻底删除type。

11.4.定制dynamic mapping

11.4.1定制dynamic策略

true: 遇到陌生字段，就进行dynamic mapping

false: 新检测到的字段将被忽略。这些字段将不会被索引，因此将无法搜索，但仍将出现在返回点击的源字段中。这些字段不会添加到映射中，必须显式添加新字段。

strict: 遇到陌生字段, 就报错

创建mapping

```
1  PUT /my_index
2  {
3      "mappings": {
4          "dynamic": "strict",
5          "properties": {
6              "title": {
7                  "type": "text"
8              },
9              "address": {
10                 "type": "object",
11                 "dynamic": "true"
12             }
13         }
14     }
15 }
```

插入数据

```
1  PUT /my_index/_doc/1
2  {
3      "title": "my article",
4      "content": "this is my article",
5      "address": {
6          "province": "guangdong",
7          "city": "guangzhou"
8      }
9  }
```

报错

```
1  {
2      "error": {
3          "root_cause": [
4              {
5                  "type": "strict_dynamic_mapping_exception",
6                  "reason": "mapping set to strict, dynamic introduction of [content]
within [_doc] is not allowed"
7              }
8          ],
9          "type": "strict_dynamic_mapping_exception",
10         "reason": "mapping set to strict, dynamic introduction of [content] within
[_doc] is not allowed"
11     },
12     "status": 400
13 }
```

#11.4.2自定义 dynamic mapping策略

es会根据传入的值, 推断类型。

JSON datatype	Elasticsearch datatype
<code>null</code>	No field is added.
<code>true</code> or <code>false</code>	<code>boolean</code> field
floating point number	<code>float</code> field
integer	<code>long</code> field
object	<code>object</code> field
array	Depends on the first non- <code>null</code> value in the array.
string	Either a <code>date</code> field (if the value passes <code>date detection</code>), a <code>double</code> or <code>long</code> field (if the value passes <code>numeric detection</code>) or a <code>text</code> field, with a <code>keyword</code> sub-field.

date_detection 日期探测

默认会按照一定格式识别date，比如yyyy-MM-dd。但是如果某个field先过来一个2017-01-01的值，就会被自动dynamic mapping成date，后面如果再来一个"hello world"之类的值，就会报错。可以手动关闭某个type的date_detection，如果有需要，自己手动指定某个field为date类型。

```

1  PUT /my_index
2  {
3    "mappings": {
4      "date_detection": false,
5      "properties": {
6        "title": {
7          "type": "text"
8        },
9        "address": {
10         "type": "object",
11         "dynamic": "true"
12       }
13     }
14   }
15 }
```

测试

```

1  PUT /my_index/_doc/1
2  {
3    "title": "my article",
4    "content": "this is my article",
5    "address": {
6      "province": "guangdong",
7      "city": "guangzhou"
8    },
9    "post_date": "2019-09-10"
10 }
```


查看映射

```
1 GET /my_index/_mapping
```

自定义日期格式

```
1 PUT my_index
2 {
3   "mappings": {
4     "dynamic_date_formats": ["MM/dd/yyyy"]
5   }
6 }
```

插入数据

```
1 PUT my_index/_doc/1
2 {
3   "create_date": "09/25/2019"
4 }
```

numeric_detection 数字探测

虽然json支持本机浮点和整数数据类型，但某些应用程序或语言有时可能将数字呈现为字符串。通常正确的解决方案是显式地映射这些字段，但是可以启用数字检测（默认情况下禁用）来自动完成这些操作。

```
1 PUT my_index
2 {
3   "mappings": {
4     "numeric_detection": true
5   }
6 }
```

```
1 PUT my_index/_doc/1
2 {
3   "my_float": "1.0",
4   "my_integer": "1"
5 }
```

11.4.3 定制自己的dynamic mapping template

```
1 PUT /my_index
2 {
3   "mappings": {
4     "dynamic_templates": [
5       {
6         "en": {
7           "match": "*_en",
8           "match_mapping_type": "string",
9           "mapping": {
10            "type": "text",
11            "analyzer": "english"
12          }
13        }
14      ]
15    }
16 }
```

```
17 }
```

插入数据

```
1 PUT /my_index/_doc/1
2 {
3   "title": "this is my first article"
4 }
5
6 PUT /my_index/_doc/2
7 {
8   "title_en": "this is my first article"
9 }
```

搜索

```
1 GET my_index/_search?q=first
2 GET my_index/_search?q=is
```

title没有匹配到任何的dynamic模板，默认就是standard分词器，不会过滤停用词，is会进入倒排索引，用is来搜索是可以搜索到的

title_en匹配到了dynamic模板，就是english分词器，会过滤停用词，is这种停用词就会被过滤掉，用is来搜索就搜索不到了

模板写法

```
1 PUT my_index
2 {
3   "mappings": {
4     "dynamic_templates": [
5       {
6         "integers": {
7           "match_mapping_type": "long",
8           "mapping": {
9             "type": "integer"
10          }
11        }
12      },
13      {
14        "strings": {
15          "match_mapping_type": "string",
16          "mapping": {
17            "type": "text",
18            "fields": {
19              "raw": {
20                "type": "keyword",
21                "ignore_above": 256
22              }
23            }
24          }
25        }
26      }
27    ]
28  }
29 }
```

模板参数

```

1  "match": "long_*",
2  "unmatch": "*_text",
3  "match_mapping_type": "string",
4  "path_match": "name.*",
5  "path_unmatch": "/*.middle",

```

```

1  "match_pattern": "regex",
2  "match": "^profit_\\d+$"

```

场景

1 结构化搜索

默认情况下，elasticsearch将字符串字段映射为带有子关键字字段的文本字段。但是，如果只对结构化内容进行索引，而对全文搜索不感兴趣，则可以仅将“字段”映射为“关键字”。请注意，这意味着为了搜索这些字段，必须搜索索引所用的完全相同的值。

```

1  {
2      "strings_as_keywords": {
3          "match_mapping_type": "string",
4          "mapping": {
5              "type": "keyword"
6          }
7      }
8  }

```

2 仅搜索

与前面的示例相反，如果您只关心字符串字段的全文搜索，并且不打算对字符串字段运行聚合、排序或精确搜索，您可以告诉弹性搜索将其仅映射为文本字段（这是5之前的默认行为）

```

1  {
2      "strings_as_text": {
3          "match_mapping_type": "string",
4          "mapping": {
5              "type": "text"
6          }
7      }
8  }

```

3 norms 不关心评分

norms是指标时间的评分因素。如果您不关心评分，例如，如果您从不按评分对文档进行排序，则可以在索引中禁用这些评分因子的存储并节省一些空间。

```

1  {
2      "strings_as_keywords": {
3          "match_mapping_type": "string",
4          "mapping": {
5              "type": "text",
6              "norms": false,
7              "fields": {
8                  "keyword": {
9                      "type": "keyword",
10                     "ignore_above": 256
11                 }
12             }
13         }
14     }
15 }

```

```
14     }
15 }
```

#11.5. 零停机重建索引

#11.5.1 零停机重建索引

场景：

一个field的设置是不能被修改的，如果要修改一个Field，那么应该重新按照新的mapping，建立一个index，然后将数据批量查询出来，重新用bulk api写入index中。

批量查询的时候，建议采用scroll api，并且采用多线程并发的方式来reindex数据，每次scroll就查询指定日期的一段数据，交给一个线程即可。

(1)一开始，依靠dynamic mapping，插入数据，但是不小心有些数据是2019-09-10这种日期格式的，所以title这种field被自动映射为了date类型，实际上它应该是string类型的

```
1  PUT /my_index/_doc/1
2  {
3    "title": "2019-09-10"
4  }
5
6  PUT /my_index/_doc/2
7  {
8    "title": "2019-09-11"
9  }
```

(2) 当后期向索引中加入string类型的title值的时候，就会报错

```
1  PUT /my_index/_doc/3
2  {
3    "title": "my first article"
4  }
```

报错

```
1  {
2    "error": {
3      "root_cause": [
4        {
5          "type": "mapper_parsing_exception",
6          "reason": "failed to parse [title]"
7        }
8      ],
9      "type": "mapper_parsing_exception",
10     "reason": "failed to parse [title]",
11     "caused_by": {
12       "type": "illegal_argument_exception",
13       "reason": "Invalid format: \"my first article\""
14     }
15   },
16   "status": 400
17 }
```

(3) 如果此时想修改title的类型，是不可能的

```

1  PUT /my_index/_mapping
2  {
3    "properties": {
4      "title": {
5        "type": "text"
6      }
7    }
8  }

```

报错

```

1  {
2    "error": {
3      "root_cause": [
4        {
5          "type": "illegal_argument_exception",
6          "reason": "mapper [title] of different type, current_type [date],
merged_type [text]"
7        }
8      ],
9      "type": "illegal_argument_exception",
10     "reason": "mapper [title] of different type, current_type [date], merged_type
[text]"
11   },
12   "status": 400
13 }

```

(4) 此时，唯一的办法，就是进行reindex，也就是说，重新建立一个索引，将旧索引的数据查询出来，再导入新索引。

(5) 如果说旧索引的名字，是old_index，新索引的名字是new_index，终端java应用，已经在使用old_index在操作了，难道还要去停止java应用，修改使用的index为new_index，才重新启动java应用吗？这个过程中，就会导致java应用停机，可用性降低。

(6) 所以说，给java应用一个别名，这个别名是指向旧索引的，java应用先用着，java应用先用prod_index alias来操作，此时实际指向的是旧的my_index

```

1  PUT /my_index/_alias/prod_index

```

1

(7) 新建一个index，调整其title的类型为string

```

1  PUT /my_index_new
2  {
3    "mappings": {
4      "properties": {
5        "title": {
6          "type": "text"
7        }
8      }
9    }
10 }

```

(8) 使用scroll api将数据批量查询出来

```

1  GET /my_index/_search?scroll=1m
2  {
3      "query": {
4          "match_all": {}
5      },
6      "size": 1
7  }

```

返回

```

1  {
2      "_scroll_id":
3      "DnF1ZXJ5VGhlbkZldGNoBQAAAAAAAAADpAFjRvbnNUWVZaVGpHdklqOV9zcFd6MncAAAAAAA6QRY0b25z
4      VF1WW1RqR3ZJajlfc3BXejJ3AAAAAAA0kIWNG9uc1RZV1pUakd2SWo5X3NwV3oydwAAAAAADpDFjRvb
5      nNUWVZaVGpHdklqOV9zcFd6MncAAAAAAA6RBY0b25zVF1WW1RqR3ZJajlfc3BXejJ3",
6      "took": 1,
7      "timed_out": false,
8      "_shards": {
9          "total": 5,
10         "successful": 5,
11         "failed": 0
12     },
13     "hits": {
14         "total": 3,
15         "max_score": null,
16         "hits": [
17             {
18                 "_index": "my_index",
19                 "_type": "my_type",
20                 "_id": "1",
21                 "_score": null,
22                 "_source": {
23                     "title": "2019-01-02"
24                 },
25                 "sort": [
26                     0
27                 ]
28             }
29         ]
30     }
31 }

```

(9) 采用bulk api将scroll查出来的一批数据，批量写入新索引

```

1  POST /_bulk
2  { "index": { "_index": "my_index_new", "_id": "1" }}
3  { "title": "2019-09-10" }

```

(10) 反复循环8~9，查询一批又一批的数据出来，采取bulk api将每一批数据批量写入新索引

(11) 将prod_index alias切换到my_index_new上去，java应用会直接通过index别名使用新的索引中的数据，java应用程序不需要停机，零提交，高可用

```
1  POST /_aliases
2  {
3      "actions": [
4          { "remove": { "index": "my_index", "alias": "prod_index" }},
5          { "add": { "index": "my_index_new", "alias": "prod_index" }}
6      ]
7  }
```

(12) 直接通过prod_index别名来查询, 是否ok

```
1  GET /prod_index/_search
```

11.5.2 生产实践：基于alias对client透明切换index

```
1  PUT /my_index_v1/_alias/my_index
```

client对my_index进行操作

reindex操作, 完成之后, 切换v1到v2

```
1  POST /_aliases
2  {
3      "actions": [
4          { "remove": { "index": "my_index_v1", "alias": "my_index" }},
5          { "add": { "index": "my_index_v2", "alias": "my_index" }}
6      ]
7  }
```

12. 中文分词器 IK分词器

12.1. IK分词器安装使用

12.1.1 中文分词器

standard 分词器, 仅适用于英文。

```
1  GET /_analyze
2  {
3      "analyzer": "standard",
4      "text": "中华人民共和国人民大会堂"
5  }
```

我们想要的效果是什么：中华人民共和国，人民大会堂

IK分词器就是目前最流行的es中文分词器

12.1.2 安装

官网：<https://github.com/medcl/elasticsearch-analysis-ik>

下载地址：<https://github.com/medcl/elasticsearch-analysis-ik/releases>

根据es版本下载相应版本包。

解压到 es/plugins/ik中。

重启es

12.1.3 ik分词器基础知识

ik_max_word: 会将文本做最细粒度的拆分，比如会将“中华人民共和国人民大会堂”拆分为“中华人民共和国，中华人民共和国，中华人民，中华，华人，人民共和国，人民大会堂，人民大会，大会堂”，会穷尽各种可能的组合；

ik_smart: 会做最粗粒度的拆分，比如会将“中华人民共和国人民大会堂”拆分为“中华人民共和国，人民大会堂”。

12.1.4 ik分词器的使用

存储时，使用ik_max_word，搜索时，使用ik_smart

```
1  PUT /my_index
2  {
3    "mappings": {
4      "properties": {
5        "text": {
6          "type": "text",
7          "analyzer": "ik_max_word",
8          "search_analyzer": "ik_smart"
9        }
10     }
11   }
12 }
```

搜索

```
1  GET /my_index/_search?q=中华人民共和国人民大会堂
```

12.2. ik配置文件

12.2.1 ik配置文件

ik配置文件地址：es/plugins/ik/config目录

IKAnalyzer.cfg.xml：用来配置自定义词库

main.dic：ik原生内置的中文词库，总共有27万多条，只要是这些单词，都会被分在一起

preposition.dic: 介词

quantifier.dic：放了一些单位相关的词，量词

suffix.dic：放了一些后缀

surname.dic：中国的姓氏

stopword.dic：英文停用词

ik原生最重要的两个配置文件

main.dic：包含了原生的中文词语，会按照这个里面的词语去分词

stopword.dic：包含了英文的停用词

停用词，stopword

a the and at but

一般，像停用词，会在分词的时候，直接被干掉，不会建立在倒排索引中

12.2.1 自定义词库

(1) 自己建立词库：每年都会涌现一些特殊的流行词，网红，蓝瘦香菇，喊麦，鬼畜，一般都不会在ik的原生词典里

自己补充自己的最新的词语，到ik的词库里面

IKAnalyzer.cfg.xml: ext_dict, 创建mydict.dic。

补充自己的词语，然后需要重启es，才能生效

(2) 自己建立停用词库：比如了，的，啥，么，我们可能并不想去建立索引，让人家搜索

custom/ext_stopword.dic，已经有了常用的中文停用词，可以补充自己的停用词，然后重启es

12.3. 使用mysql热更新词库

12.3.1热更新

每次都是在es的扩展词典中，手动添加新词语，很坑

(1) 每次添加完，都要重启es才能生效，非常麻烦

(2) es是分布式的，可能有数百个节点，你不能每次都一个一个节点上面去修改

es不停机，直接我们在外部某个地方添加新的词语，es中立即热加载到这些新词语

热更新的方案

(1) 基于ik分词器原生支持的热更新方案，部署一个web服务器，提供一个http接口，通过modified和tag两个http响应头，来提供词语的热更新

(2) 修改ik分词器源码，然后手动支持从mysql中每隔一定时间，自动加载新的词库

用第二种方案，第一种，ik git社区官方都不建议采用，觉得不太稳定

12.3.2步骤

1、下载源码

<https://github.com/medcl/elasticsearch-analysis-ik/releases>

ik分词器，是个标准的java maven工程，直接导入eclipse就可以看到源码

2、修改源

org.wltea.analyzer.dic.Dictionary类，160行Dictionary单例类的初始化方法，在这里需要创建一个我们自定义的线程，并且启动它

org.wltea.analyzer.dic.HotDictReloadThread类：就是死循环，不断调用Dictionary.getSingleton().reloadMainDict()，去重新加载词典

Dictionary类，399行：this.loadMySQLExtDict(); 加载mysql字典。

Dictionary类，609行：this.loadMySQLStopwordDict();加载mysql停用词

config下jdbc-reload.properties。mysql配置文件

3、mvn package打包代码

target\releases\elasticsearch-analysis-ik-7.3.0.zip

4、解压缩ik压缩包

将mysql驱动jar，放入ik的目录下

5、修改jdbc相关配置

6、重启es

观察日志，日志中就会显示我们打印的那些东西，比如加载了什么配置，加载了什么词语，什么停用词

7、在mysql中添加词库与停用词

8、分词实验，验证热更新生效

```
1 GET /_analyze
2 {
3   "analyzer": "ik_smart",
4   "text": "传智播客"
5 }
```

#13. java api 实现索引管理

代码

```
1 package com.ydlclass.es;
2
3 import org.elasticsearch.action.ActionListener;
4 import org.elasticsearch.action.admin.indices.alias.Alias;
5 import org.elasticsearch.action.admin.indices.close.CloseIndexRequest;
6 import org.elasticsearch.action.admin.indices.delete.DeleteIndexRequest;
7 import org.elasticsearch.action.admin.indices.open.OpenIndexRequest;
8 import org.elasticsearch.action.admin.indices.open.OpenIndexResponse;
9 import org.elasticsearch.action.support.ActiveShardCount;
10 import org.elasticsearch.action.support.master.AcknowledgedResponse;
11 import org.elasticsearch.client.IndicesClient;
12 import org.elasticsearch.client.RequestOptions;
13 import org.elasticsearch.client.RestHighLevelClient;
14 import org.elasticsearch.client.indices.CreateIndexRequest;
15 import org.elasticsearch.client.indices.CreateIndexResponse;
16 import org.elasticsearch.client.indices.GetIndexRequest;
17 import org.elasticsearch.common.settings.Settings;
18 import org.elasticsearch.common.unit.TimeValue;
19 import org.elasticsearch.common.xcontent.XContentType;
20 import org.junit.Test;
21 import org.junit.runner.RunWith;
22 import org.springframework.beans.factory.annotation.Autowired;
23 import org.springframework.boot.test.context.SpringBootTest;
24 import org.springframework.test.context.junit4.SpringRunner;
25
26 import java.io.IOException;
27
28 /**
29
30  - @author Administrator
31
```

```

32 - @version 1.0
33   **/
34   @SpringBootTest
35   @RunWith(SpringRunner.class)
36   public class TestIndex {
37
38       @Autowired
39       RestHighLevelClient client;
40
41       // @Autowired
42       // RestClient restClient;
43
44       ...
45       //创建索引
46       @Test
47       public void testCreateIndex() throws IOException {
48           //创建索引对象
49           CreateIndexRequest createIndexRequest = new
CreateIndexRequest("ydlclass_book");
50           //设置参数
51           createIndexRequest.settings(Settings.builder().put("number_of_shards",
"1").put("number_of_replicas", "0"));
52           //指定映射1
53           createIndexRequest.mapping(" {\n" +
54               " \t\"properties\": {\n" +
55               "         \"name\":{\n" +
56               "         \"type\": \"keyword\"\n" +
57               "     },\n" +
58               "     \"description\": {\n" +
59               "         \"type\": \"text\"\n" +
60               "     },\n" +
61               "     \"price\":{\n" +
62               "         \"type\": \"long\"\n" +
63               "     },\n" +
64               "     \"pic\":{\n" +
65               "         \"type\": \"text\", \n" +
66               "         \"index\": false\n" +
67               "     }\n" +
68               " \t}\n" +
69               "}", XContentType.JSON);
70
71           //指定映射2
72           ...
73
74           // Map<String, Object> message = new HashMap<>();
75           // message.put("type", "text");
76           // Map<String, Object> properties = new HashMap<>();
77           // properties.put("message", message);
78           // Map<String, Object> mapping = new HashMap<>();
79           // mapping.put("properties", properties);
80           // createIndexRequest.mapping(mapping);
81
82           ...
83           //指定映射3
84           ...
85
86           // XContentBuilder builder = XContentFactory.jsonBuilder();
87           // builder.startObject();

```

```

88     {
89         builder.startObject("properties");
90         {
91             builder.startObject("message");
92             {
93                 builder.field("type", "text");
94             }
95             builder.endObject();
96         }
97         builder.endObject();
98     }
99     builder.endObject();
100     createIndexRequest.mapping(builder);
101
102     ...
103     //设置别名
104     createIndexRequest.alias(new Alias("ydlclass_index_new"));
105
106     // 额外参数
107     //设置超时时间
108     createIndexRequest.setTimeout(TimeValue.timeValueMinutes(2));
109     //设置主节点超时时间
110     createIndexRequest.setMasterTimeout(TimeValue.timeValueMinutes(1));
111     //在创建索引API返回响应之前等待的活动分片副本的数量，以int形式表示
112     createIndexRequest.waitForActiveShards(ActiveShardCount.from(2));
113     createIndexRequest.waitForActiveShards(ActiveShardCount.DEFAULT);
114
115     //操作索引的客户端
116     IndicesClient indices = client.indices();
117     //执行创建索引库
118     CreateIndexResponse createIndexResponse = indices.create(createIndexRequest,
RequestOptions.DEFAULT);
119
120     //得到响应（全部）
121     boolean acknowledged = createIndexResponse.isAcknowledged();
122     //得到响应 指示是否在超时前为索引中的每个分片启动了所需数量的碎片副本
123     boolean shardsAcknowledged = createIndexResponse.isShardsAcknowledged();
124
125     System.out.println("!!!!!!!!!!!!!!!!!!!!!!!!!!!!!!" + acknowledged);
126     System.out.println(shardsAcknowledged);
127
128 }
129
130 //异步新增索引
131 @Test
132 public void testCreateIndexAsync() throws IOException {
133     //创建索引对象
134     CreateIndexRequest createIndexRequest = new
CreateIndexRequest("ydlclass_book2");
135     //设置参数
136     createIndexRequest.settings(Settings.builder().put("number_of_shards",
"1").put("number_of_replicas", "0"));
137     //指定映射1
138     createIndexRequest.mapping(" {\n" +
139         " \t\"properties\": {\n" +
140         "         \"name\":{\n" +
141         "             \"type\": \"keyword\"\n" +
142         "         },\n" +

```

```

143         "description\": {\n" +
144         "            \"type\": \"text\"\n" +
145         "        },\n" +
146         "        \"price\":{\n" +
147         "            \"type\": \"long\"\n" +
148         "        },\n" +
149         "        \"pic\":{\n" +
150         "            \"type\": \"text\", \n" +
151         "            \"index\": false\n" +
152         "        }\n" +
153         "    }\n" +
154     "}", XContentType.JSON);
155
156     //监听方法
157     ActionListener<CreateIndexResponse> listener =
158         new ActionListener<CreateIndexResponse>() {
159
160         @Override
161         public void onResponse(CreateIndexResponse createIndexResponse)
162     {
163         System.out.println("!!!!!!!创建索引成功");
164         System.out.println(createIndexResponse.toString());
165     }
166
167     @Override
168     public void onFailure(Exception e) {
169         System.out.println("!!!!!!!创建索引失败");
170         e.printStackTrace();
171     }
172     };
173
174     //操作索引的客户端
175     IndicesClient indices = client.indices();
176     //执行创建索引库
177     indices.createAsync(createIndexRequest, RequestOptions.DEFAULT, listener);
178
179     try {
180         Thread.sleep(5000);
181     } catch (InterruptedException e) {
182         e.printStackTrace();
183     }
184     ...
185     ...
186     }
187     ...
188     ...
189     ...
190     //删除索引库
191     @Test
192     public void testDeleteIndex() throws IOException {
193         //删除索引对象
194         DeleteIndexRequest deleteIndexRequest = new
195         DeleteIndexRequest("ydlclass_book2");
196         //操作索引的客户端
197         IndicesClient indices = client.indices();
198         //执行删除索引

```

```

198     AcknowledgedResponse delete = indices.delete(deleteIndexRequest,
RequestOptions.DEFAULT);
199     //得到响应
200     boolean acknowledged = delete.isAcknowledged();
201     System.out.println(acknowledged);
202
203 }
204
205 //异步删除索引库
206 @Test
207 public void testDeleteIndexAsync() throws IOException {
208     //删除索引对象
209     DeleteIndexRequest deleteIndexRequest = new
DeleteIndexRequest("ydlclass_book2");
210     //操作索引的客户端
211     IndicesClient indices = client.indices();
212
213     //监听方法
214     ActionListener<AcknowledgedResponse> listener =
215         new ActionListener<AcknowledgedResponse>() {
216             @Override
217             public void onResponse(AcknowledgedResponse deleteIndexResponse)
218             {
219                 System.out.println("!!!!!!!删除索引成功");
220                 System.out.println(deleteIndexResponse.toString());
221             }
222             @Override
223             public void onFailure(Exception e) {
224                 System.out.println("!!!!!!!删除索引失败");
225                 e.printStackTrace();
226             }
227         };
228     //执行删除索引
229     indices.deleteAsync(deleteIndexRequest, RequestOptions.DEFAULT, listener);
230
231     try {
232         Thread.sleep(5000);
233     } catch (InterruptedException e) {
234         e.printStackTrace();
235     }
236
237 }
238
239 // Indices Exists API
240 @Test
241 public void testExistIndex() throws IOException {
242     GetIndexRequest request = new GetIndexRequest("ydlclass_book");
243     request.local(false); //从主节点返回本地信息或检索状态
244     request.humanReadable(true); //以适合人类的格式返回结果
245     request.includeDefaults(false); //是否返回每个索引的所有默认设置
246
247     boolean exists = client.indices().exists(request, RequestOptions.DEFAULT);
248     System.out.println(exists);
249 }
250 ...
251
252 ...

```

```

253 // Indices Open API
254 @Test
255 public void testOpenIndex() throws IOException {
256     OpenIndexRequest request = new OpenIndexRequest("ydlclass_book");
257
258     OpenIndexResponse openIndexResponse = client.indices().open(request,
RequestOptions.DEFAULT);
259     boolean acknowledged = openIndexResponse.isAcknowledged();
260     System.out.println("!!!!!!!!!!"+acknowledged);
261 }
262
263 // Indices Close API
264 @Test
265 public void testCloseIndex() throws IOException {
266     CloseIndexRequest request = new CloseIndexRequest("index");
267     AcknowledgedResponse closeIndexResponse = client.indices().close(request,
RequestOptions.DEFAULT);
268     boolean acknowledged = closeIndexResponse.isAcknowledged();
269     System.out.println("!!!!!!!!!!"+acknowledged);
270
271 }
272 }

```

#14. search搜索入门

#14.1. 搜索语法入门

#14.1.1 query string search

无条件搜索所有

```
1 GET /book/_search
```

```

1  {
2    "took" : 969,
3    "timed_out" : false,
4    "_shards" : {
5      "total" : 1,
6      "successful" : 1,
7      "skipped" : 0,
8      "failed" : 0
9    },
10   "hits" : {
11     "total" : {
12       "value" : 3,
13       "relation" : "eq"
14     },
15     "max_score" : 1.0,
16     "hits" : [
17       {
18         "_index" : "book",
19         "_type" : "_doc",
20         "_id" : "1",
21         "_score" : 1.0,

```

```

22         "_source" : {
23             "name" : "Bootstrap开发",
24             "description" : "Bootstrap是由Twitter推出的一个前台页面开发css框架，是一个非常流行的开发框架，此框架集成了多种页面效果。此开发框架包含了大量的CSS、JS程序代码，可以帮助开发者（尤其是不擅长css页面开发的程序人员）轻松的实现一个css，不受浏览器限制的精美界面css效果。",
25             "studymodel" : "201002",
26             "price" : 38.6,
27             "timestamp" : "2019-08-25 19:11:35",
28             "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
29             "tags" : [
30                 "bootstrap",
31                 "dev"
32             ]
33         },
34     },
35     {
36         "_index" : "book",
37         "_type" : "_doc",
38         "_id" : "2",
39         "_score" : 1.0,
40         "_source" : {
41             "name" : "java编程思想",
42             "description" : "java语言是世界第一编程语言，在软件开发领域使用人数最多。",
43             "studymodel" : "201001",
44             "price" : 68.6,
45             "timestamp" : "2019-08-25 19:11:35",
46             "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
47             "tags" : [
48                 "java",
49                 "dev"
50             ]
51         },
52     },
53     {
54         "_index" : "book",
55         "_type" : "_doc",
56         "_id" : "3",
57         "_score" : 1.0,
58         "_source" : {
59             "name" : "spring开发基础",
60             "description" : "spring 在java领域非常流行，java程序员都在用。",
61             "studymodel" : "201001",
62             "price" : 88.6,
63             "timestamp" : "2019-08-24 19:11:35",
64             "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
65             "tags" : [
66                 "spring",
67                 "java"
68             ]
69         },
70     }
71 ]
72 }
73 }

```

解释

took: 耗费了几毫秒

timed_out: 是否超时, 这里是没有

_shards: 到几个分片搜索, 成功几个, 跳过几个, 失败几个。

hits.total: 查询结果的数量, 3个document

hits.max_score: score的含义, 就是document对于一个search的相关度的匹配分数, 越相关, 就越匹配, 分数也高

hits.hits: 包含了匹配搜索的document的所有详细数据

14.1.2 传参

与http请求传参类似

```
1 GET /book/_search?q=name:java&sort=price:desc
```

类比sql: select * from book where name like ' %java%' order by price desc

```
1  {
2    "took" : 2,
3    "timed_out" : false,
4    "_shards" : {
5      "total" : 1,
6      "successful" : 1,
7      "skipped" : 0,
8      "failed" : 0
9    },
10   "hits" : {
11     "total" : {
12       "value" : 1,
13       "relation" : "eq"
14     },
15     "max_score" : null,
16     "hits" : [
17       {
18         "_index" : "book",
19         "_type" : "_doc",
20         "_id" : "2",
21         "_score" : null,
22         "_source" : {
23           "name" : "java编程思想",
24           "description" : "java语言是世界第一编程语言, 在软件开发领域使用人数最多。",
25           "studymodel" : "201001",
26           "price" : 68.6,
27           "timestamp" : "2019-08-25 19:11:35",
28           "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
29           "tags" : [
30             "java",
31             "dev"
32           ]
33         },
34         "sort" : [
35           68.6
36         ]
37       }
38     ]
39   }
40 }
```

14.1.3 图解timeout

GET /book/_search?timeout=10ms

全局设置：配置文件中设置 search.default_search_timeout：100ms。默认不超时。

14.2. multi-index 多索引搜索

14.2.1 multi-index 搜索模式

告诉你如何一次性搜索多个index和多个type下的数据

- 1 /_search: 所有索引下的所有数据都搜索出来
- 2 /index1/_search: 指定一个index，搜索其下所有的数据
- 3 /index1,index2/_search: 同时搜索两个index下的数据
- 4 /index*/_search: 按照通配符去匹配多个索引

应用场景：生产环境log索引可以按照日期分开。

log_to_es_20190910

log_to_es_20190911

log_to_es_20180910

14.2.2 初步图解一下简单的搜索原理

搜索原理初步图解

14.3. 分页搜索

14.3.1 分页搜索的语法

sql: select * from book limit 1,5

size, from

GET /book/_search?size=10

GET /book/_search?size=10&from=0

GET /book/_search?size=10&from=20

GET /book_search?from=0&size=3

14.3.2 deep paging

什么是deep paging

根据相关度评分倒排序，所以分页过深，协调节点会将大量数据聚合分析。

deep paging 性能问题

1消耗网络带宽，因为所搜过深的话，各 shard 要把数据传递给 coordinate node，这个过程是有大量数据传递的，消耗网络。

2消耗内存，各 shard 要把数据传送给 coordinate node，这个传递回来的数据，是被 coordinate node 保存在内存中的，这样会大量消耗内存。

3消耗cup，coordinate node 要把传回来的数据进行排序，这个排序过程很消耗cpu。所以：鉴于deep paging的性能问题，所有应尽量减少使用。

14.4. query string基础语法

14.4.1query string基础语法

GET /book/_search?q=name:java

GET /book/_search?q=+name:java

GET /book/_search?q=-name:java

一个是掌握q=field:search content的语法，还有一个是掌握+和-的含义

14.4.2、_all metadata的原理和作用

```
1 GET /book/_search?q=java
```

直接可以搜索所有的field，任意一个field包含指定的关键字就可以搜索出来。我们在进行中搜索的时候，难道是对document中的每一个field都进行一次搜索吗？不是的。

es中 *all*元数据。建立索引的时候，插入一条 *document*，*es*会将所有的 *field*值经行全量分词，把这些分词，放到 *all field*中。在搜索的时候，没有指定 *field*，就在 *_all*搜索。

举例

```
1 {
2     name:jack
3     email:123@qq.com
4     address:beijing
5 }
```

_all : jack, 123@qq.com ,beijing

14.5. query DSL入门

14.5.1 DSL

query string 后边的参数原来越多，搜索条件越来越复杂，不能满足需求。

GET /book/_search?q=name:java&size=10&from=0&sort=price:desc

DSL:Domain Specified Language，特定领域的语言

es特有的搜索语言，可在请求体中携带搜索条件，功能强大。

查询全部 GET /book/_search

```
1  GET /book/_search
2  {
3    "query": { "match_all": {} }
4  }
```

排序 GET /book/_search?sort=price:desc

```
1  GET /book/_search
2  {
3    "query" : {
4      "match" : {
5        "name" : " java"
6      }
7    },
8    "sort": [
9      { "price": "desc" }
10   ]
11 }
```

分页查询 GET /book/_search?size=10&from=0

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

指定返回字段 GET /book/_search?_source=name,studymodel

```
1  GET /book/_search
2  {
3    "query": { "match_all": {} },
4    "_source": ["name", "studymodel"]
5  }
```

通过组合以上各种类型查询，实现复杂查询。

#14.5.2. Query DSL语法

```
1  {
2    QUERY_NAME: {
3      ARGUMENT: VALUE,
4      ARGUMENT: VALUE, ...
5    }
6  }
```

```
1  {
2    QUERY_NAME: {
3      FIELD_NAME: {
4        ARGUMENT: VALUE,
5        ARGUMENT: VALUE, ...
6      }
7    }
8  }
```

```
1 GET /test_index/_search
2 {
3   "query": {
4     "match": {
5       "test_field": "test"
6     }
7   }
8 }
```

#14.5.3 组合多个搜索条件

搜索需求: title必须包含elasticsearch, content可以包含elasticsearch也可以不包含, author_id必须不为111

sql where and or !=

初始数据:

```
1 POST /website/_doc/1
2 {
3   "title": "my hadoop article",
4   "content": "hadoop is very bad",
5   "author_id": 111
6 }
7
8 POST /website/_doc/2
9 {
10   "title": "my elasticsearch article",
11   "content": "es is very bad",
12   "author_id": 112
13 }
14 POST /website/_doc/3
15 {
16   "title": "my elasticsearch article",
17   "content": "es is very goods",
18   "author_id": 111
19 }
```

搜索:

```
1 GET /website/_doc/_search
2 {
3   "query": {
4     "bool": {
5       "must": [
6         {
7           "match": {
8             "title": "elasticsearch"
9           }
10        }
11      ],
12       "should": [
13         {
14           "match": {
15             "content": "elasticsearch"
16           }
17         }
18       ],

```

```

19     "must_not": [
20     {
21         "match": {
22             "author_id": 111
23         }
24     }
25 ]
26 }
27 }
28 }

```

返回:

```

1  {
2    "took" : 488,
3    "timed_out" : false,
4    "_shards" : {
5      "total" : 1,
6      "successful" : 1,
7      "skipped" : 0,
8      "failed" : 0
9    },
10   "hits" : {
11     "total" : {
12       "value" : 1,
13       "relation" : "eq"
14     },
15     "max_score" : 0.47000363,
16     "hits" : [
17       {
18         "_index" : "website",
19         "_type" : "_doc",
20         "_id" : "2",
21         "_score" : 0.47000363,
22         "_source" : {
23           "title" : "my elasticsearch article",
24           "content" : "es is very bad",
25           "author_id" : 112
26         }
27       }
28     ]
29   }
30 }

```

更复杂的搜索需求:

select * from test_index where name='tom' or (hired =true and (personality ='good' and rude != true))

```

1  GET /test_index/_search
2  {
3    "query": {
4      "bool": {
5        "must": { "match":{ "name": "tom" }},
6        "should": [
7          { "match":{ "hired": true }},
8          { "bool": {
9            "must":{ "match": { "personality": "good" }},
10           "must_not": { "match": { "rude": true }}

```

```

11         }}
12     ],
13     "minimum_should_match": 1
14 }
15 }
16 }

```

#14.6. full-text search 全文检索

#14.6.1 全文检索

重新创建book索引

```

1  PUT /book/
2  {
3    "settings": {
4      "number_of_shards": 1,
5      "number_of_replicas": 0
6    },
7    "mappings": {
8      "properties": {
9        "name": {
10         "type": "text",
11         "analyzer": "ik_max_word",
12         "search_analyzer": "ik_smart"
13       },
14       "description": {
15         "type": "text",
16         "analyzer": "ik_max_word",
17         "search_analyzer": "ik_smart"
18       },
19       "studymodel": {
20         "type": "keyword"
21       },
22       "price": {
23         "type": "double"
24       },
25       "timestamp": {
26         "type": "date",
27         "format": "yyyy-MM-dd HH:mm:ss||yyyy-MM-dd||epoch_millis"
28       },
29       "pic": {
30         "type": "text",
31         "index": false
32       }
33     }
34   }
35 }

```

插入数据

```

1  PUT /book/_doc/1
2  {
3    "name": "Bootstrap开发",

```

```

4  "description": "Bootstrap是由Twitter推出的一个前台页面开发css框架，是一个非常流行的开发框
   架，此框架集成了多种页面效果。此开发框架包含了大量的CSS、JS程序代码，可以帮助开发者（尤其是不擅长
   css页面开发的程序人员）轻松的实现一个css，不受浏览器限制的精美界面css效果。",
5  "studymodel": "201002",
6  "price":38.6,
7  "timestamp":"2019-08-25 19:11:35",
8  "pic":"group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
9  "tags": [ "bootstrap", "dev"]
10 }
11
12 PUT /book/_doc/2
13 {
14  "name": "java编程思想",
15  "description": "java语言是世界第一编程语言，在软件开发领域使用人数最多。",
16  "studymodel": "201001",
17  "price":68.6,
18  "timestamp":"2019-08-25 19:11:35",
19  "pic":"group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
20  "tags": [ "java", "dev"]
21 }
22
23 PUT /book/_doc/3
24 {
25  "name": "spring开发基础",
26  "description": "spring 在java领域非常流行，java程序员都在用。",
27  "studymodel": "201001",
28  "price":88.6,
29  "timestamp":"2019-08-24 19:11:35",
30  "pic":"group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
31  "tags": [ "spring", "java"]
32 }

```

搜索

```

1  GET  /book/_search
2  {
3      "query" : {
4          "match" : {
5              "description" : "java程序员"
6          }
7      }
8  }

```

#14.6.2 _score初探

```

1  {
2      "took" : 1,
3      "timed_out" : false,
4      "_shards" : {
5          "total" : 1,
6          "successful" : 1,
7          "skipped" : 0,
8          "failed" : 0
9      },
10     "hits" : {
11         "total" : {
12             "value" : 2,

```



```

13     "relation" : "eq"
14 },
15     "max_score" : 2.137549,
16     "hits" : [
17     {
18         "_index" : "book",
19         "_type" : "_doc",
20         "_id" : "3",
21         "_score" : 2.137549,
22         "_source" : {
23             "name" : "spring开发基础",
24             "description" : "spring 在java领域非常流行, java程序员都在用。",
25             "studymodel" : "201001",
26             "price" : 88.6,
27             "timestamp" : "2019-08-24 19:11:35",
28             "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
29             "tags" : [
30                 "spring",
31                 "java"
32             ]
33         }
34     },
35     {
36         "_index" : "book",
37         "_type" : "_doc",
38         "_id" : "2",
39         "_score" : 0.57961315,
40         "_source" : {
41             "name" : "java编程思想",
42             "description" : "java语言是世界第一编程语言, 在软件开发领域使用人数最多。",
43             "studymodel" : "201001",
44             "price" : 68.6,
45             "timestamp" : "2019-08-25 19:11:35",
46             "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
47             "tags" : [
48                 "java",
49                 "dev"
50             ]
51         }
52     }
53 ]
54 }
55 }

```

结果分析

1、建立索引时, description字段 term倒排索引

java 2,3

程序员 3

2、搜索时, 直接找description中含有java的文档 2,3, 并且3号文档含有两个java字段, 一个程序员, 所以得分高, 排在前面。2号文档含有一个java, 排在后面。

#14.7. DSL 语法练习

#14.7.1 match_all

```
1  GET /book/_search
2  {
3      "query": {
4          "match_all": {}
5      }
6  }
```

#14.7.2 match

```
1  GET /book/_search
2  {
3      "query": {
4          "match": {
5              "description": "java程序员"
6          }
7      }
8  }
```

#14.7.3 multi_match

```
1  GET /book/_search
2  {
3      "query": {
4          "multi_match": {
5              "query": "java程序员",
6              "fields": ["name", "description"]
7          }
8      }
9  }
```

#14.7.4、range query 范围查询

```
1  GET /book/_search
2  {
3      "query": {
4          "range": {
5              "price": {
6                  "gte": 80,
7                  "lte": 90
8              }
9          }
10     }
11 }
```

#14.7.5、term query

字段为keyword时，存储和搜索都不分词

```
1 GET /book/_search
2 {
3   "query": {
4     "term": {
5       "description": "java程序员"
6     }
7   }
8 }
```

#14.7.6、terms query

```
1 GET /book/_search
2 {
3   "query": { "terms": { "tag": [ "search", "full_text", "nosql" ] }}
4 }
```

#14.7.7、exist query 查询有某些字段值的文档

```
1 GET /_search
2 {
3   "query": {
4     "exists": {
5       "field": "join_date"
6     }
7   }
8 }
```

#14.7.8、Fuzzy query

返回包含与搜索词类似的词的文档，该词由Levenshtein编辑距离度量。

包括以下几种情况：

- 更改角色 (box→fox)
- 删除字符 (aple→apple)
- 插入字符 (sick→sic)
- 调换两个相邻字符 (ACT→CAT)

```
1 GET /book/_search
2 {
3   "query": {
4     "fuzzy": {
5       "description": {
6         "value": "jave"
7       }
8     }
9   }
10 }
```

#14.7.9、IDs

```
1  GET /book/_search
2  {
3      "query": {
4          "ids" : {
5              "values" : ["1", "4", "100"]
6          }
7      }
8  }
```

#14.7.10、prefix 前缀查询

```
1  GET /book/_search
2  {
3      "query": {
4          "prefix": {
5              "description": {
6                  "value": "spring"
7              }
8          }
9      }
10 }
```

#14.7.11、regexp query 正则查询

```
1  GET /book/_search
2  {
3      "query": {
4          "regexp": {
5              "description": {
6                  "value": "j.*a",
7                  "flags" : "ALL",
8                  "max_determinized_states": 10000,
9                  "rewrite": "constant_score"
10             }
11         }
12     }
13 }
```

#14.8. Filter

#14.8.1 filter与query示例

需求：用户查询description中有"java程序员"，并且价格大于80小于90的数据。

```
1  GET /book/_search
2  {
3      "query": {
4          "bool": {
5              "must": [
6                  {
7                      "match": {
8                          "description": "java程序员"
9                      }
10                 }
11             ]
12         }
13     }
```

```

10         },
11         {
12             "range": {
13                 "price": {
14                     "gte": 80,
15                     "lte": 90
16                 }
17             }
18         }
19     ]
20 }
21 }
22 }

```

使用filter:

```

1  GET /book/_search
2  {
3      "query": {
4          "bool": {
5              "must": [
6                  {
7                      "match": {
8                          "description": "java程序员"
9                      }
10                 }
11             ],
12             "filter": {
13                 "range": {
14                     "price": {
15                         "gte": 80,
16                         "lte": 90
17                     }
18                 }
19             }
20         }
21     }
22 }

```

14.8.2 filter与query对比

filter，仅仅只是按照搜索条件过滤出需要的数据而已，不计算任何相关度分数，对相关度没有任何影响。

query，会去计算每个document相对于搜索条件的相关度，并按照相关度进行排序。

应用场景：

一般来说，如果你是在进行搜索，需要将最匹配搜索条件的数据先返回，那么用query 如果你只是要根据一些条件筛选出一部分数据，不关注其排序，那么用filter

14.8.3 filter与query性能

filter，不需要计算相关度分数，不需要按照相关度分数进行排序，同时还有内置的自动cache最常使用filter的数据

query，相反，要计算相关度分数，按照分数进行排序，而且无法cache结果

#14.9. 定位错误语法

验证错误语句：

```
1 GET /book/_validate/query?explain
2 {
3   "query": {
4     "mach": {
5       "description": "java程序员"
6     }
7   }
8 }
```

返回：

```
1 {
2   "valid" : false,
3   "error" : "org.elasticsearch.common.ParsingException: no [query] registered for
4     [mach]"
5 }
```

正确

```
1 GET /book/_validate/query?explain
2 {
3   "query": {
4     "match": {
5       "description": "java程序员"
6     }
7   }
8 }
```

返回

```
1 {
2   "_shards" : {
3     "total" : 1,
4     "successful" : 1,
5     "failed" : 0
6   },
7   "valid" : true,
8   "explanations" : [
9     {
10      "index" : "book",
11      "valid" : true,
12      "explanation" : "description:java description:程序员"
13    }
14  ]
15 }
```

一般用在那种特别复杂庞大的搜索下，比如你一下子写了上百行的搜索，这个时候可以先用validate api去验证一下，搜索是否合法。

合法以后，explain就像mysql的执行计划，可以看到搜索的目标等信息。

14.10. 定制排序规则

14.10.1 默认排序规则

默认情况下，是按照_score降序排序的

然而，某些情况下，可能没有有用的_score，比如说filter

```
1  GET book/_search
2  {
3    "query": {
4      "bool": {
5        "must": [
6          {
7            "match": {
8              "description": "java程序员"
9            }
10         }
11       ]
12     }
13   }
14 }
```

当然，也可以是constant_score

14.10.2 定制排序规则

相当于sql中order by ?sort=sprice:desc

```
1  GET /book/_search
2  {
3    "query": {
4      "constant_score": {
5        "filter": {
6          "term": {
7            "studymodel": "201001"
8          }
9        }
10     },
11     "sort": [
12       {
13         "price": {
14           "order": "asc"
15         }
16       }
17     ]
18   }
19 }
```

14.11. Text字段排序问题

如果对一个text field进行排序，结果往往不准确，因为分词后是多个单词，再排序就不是我们想要的结果了。

通常解决方案是，将一个text field建立两次索引，一个分词，用来进行搜索；一个不分词，用来进行排序。

fielddate:true

```
1  PUT /website
2  {
3    "mappings": {
4      "properties": {
5        "title": {
6          "type": "text",
7          "fields": {
8            "keyword": {
9              "type": "keyword"
10           }
11        }
12      },
13      "content": {
14        "type": "text"
15      },
16      "post_date": {
17        "type": "date"
18      },
19      "author_id": {
20        "type": "long"
21      }
22    }
23  }
24 }
```

插入数据

```
1  PUT /website/_doc/1
2  {
3    "title": "first article",
4    "content": "this is my second article",
5    "post_date": "2019-01-01",
6    "author_id": 110
7  }
8
9  PUT /website/_doc/2
10 {
11   "title": "second article",
12   "content": "this is my second article",
13   "post_date": "2019-01-01",
14   "author_id": 110
15 }
16
17 PUT /website/_doc/3
18 {
19   "title": "third article",
20   "content": "this is my third article",
21   "post_date": "2019-01-02",
22   "author_id": 110
23 }
```

搜索


```

1  GET /website/_search
2  {
3    "query": {
4      "match_all": {}
5    },
6    "sort": [
7      {
8        "title.keyword": {
9          "order": "desc"
10       }
11     ]
12   }
13 }

```

#14.12. Scroll分批查询

场景：下载某一个索引中1亿条数据，到文件或是数据库。

不能一下全查出来，系统内存溢出。所以使用scroll滚动搜索技术，一批一批查询。

scroll搜索会在第一次搜索的时候，保存一个当时的视图快照，之后只会基于该旧的视图快照提供数据搜索，如果这个期间数据变更，是不会让用户看到的

每次发送scroll请求，我们还需要指定一个scroll参数，指定一个时间窗口，每次搜索请求只要在这个时间窗口内能完成就可以了。

搜索

```

1  GET /book/_search?scroll=1m
2  {
3    "query": {
4      "match_all": {}
5    },
6    "size": 3
7  }

```

返回

```

1  {
2    "_scroll_id" :
3    "DXF1ZXJ5QW5kRmV0Y2gBAAAAAAAAAAMOKWTURBNDUtcjZTVUdKMFP5cX1oVE10QQ==",
4    "took" : 3,
5    "timed_out" : false,
6    "_shards" : {
7      "total" : 1,
8      "successful" : 1,
9      "skipped" : 0,
10     "failed" : 0
11   },
12   "hits" : {
13     "total" : {
14       "value" : 3,
15       "relation" : "eq"
16     },
17     "max_score" : 1.0,
18     "hits" : [
19

```

```
20     }
21 }
```

获得的结果会有一个scroll_id，下一次再发送scroll请求的时候，必须带上这个scroll_id

```
1  GET /_search/scroll
2  {
3      "scroll": "1m",
4      "scroll_id" :
5      "DXF1ZXJ5QW5kRmV0Y2gBAAAAAAAAAMOkWTURBNdUtcjZTVUdKMFP5cXloVEl0QQ=="
6  }
```

与分页区别：

分页给用户看的 deep paging

scroll是用户系统内部操作，如下载批量数据，数据转移。零停机改变索引映射。

#15. java api实现搜索

#16. 评分机制详解

#16.1. 评分机制 TF\IDF

#16.1.1 算法介绍

relevance score算法，简单来说，就是计算出，一个索引中的文本，与搜索文本，他们之间的关联匹配程度。

Elasticsearch使用的是 term frequency/inverse document frequency算法，简称为TF/IDF算法。TF词频(Term Frequency)，IDF逆向文件频率(Inverse Document Frequency)

Term frequency：搜索文本中的各个词条在field文本中出现了多少次，出现次数越多，就越相关。

$$\text{公式：} \quad tf_{ij} = \frac{n_{i,j}}{\sum_k n_{k,j}} \quad \text{即：} \quad TF_w = \frac{\text{在某一类中词条}w\text{出现的次数}}{\text{该类中所有的词条数目}}$$

举例：搜索请求：hello world

doc1 : hello you and me,and world is very good.

doc2 : hello,how are you

Inverse document frequency：搜索文本中的各个词条在整个索引的所有文档中出现了多少次，出现的次数越多，就越不相关。

$$\text{公式：} \quad idf_i = \log \frac{|D|}{|\{j : t_i \in d_j\}|}$$

$$IDF = \log\left(\frac{\text{语料库的文档总数}}{\text{包含词条}w\text{的文档数} + 1}\right)$$
,分母之所以要加1,是为了避免分母为0
即: <https://blog.csdn.net/asialeebird>

举例: 搜索请求: hello world

doc1 : hello ,today is very good

doc2 : hi world ,how are you

整个index中1亿条数据。hello的document 1000个, 有world的document 有100个。

doc2 更相关

Field-length norm: field长度, field越长, 相关度越弱

举例: 搜索请求: hello world

doc1 : {"title":"hello article","content ":"balabalabal 1万个"}

doc2 : {"title":"my article","content ":"balabalabal 1万个,world"}

16.1.2 _score是如何被计算出来的

```
1 GET /book/_search?explain=true
2 {
3   "query": {
4     "match": {
5       "description": "java程序员"
6     }
7   }
8 }
```

返回

```
1 {
2   "took" : 5,
3   "timed_out" : false,
4   "_shards" : {
5     "total" : 1,
6     "successful" : 1,
7     "skipped" : 0,
8     "failed" : 0
9   },
10  "hits" : {
11    "total" : {
12      "value" : 2,
13      "relation" : "eq"
14    },
15    "max_score" : 2.137549,
16    "hits" : [
17      {
18        "_shard" : "[book][0]",
19        "_node" : "MDA45-r6SUGJ0ZyqyhTINA",
20        "_index" : "book",
21        "_type" : "_doc",
22        "_id" : "3",
23        "_score" : 2.137549,
```

```

24         "_source" : {
25             "name" : "spring开发基础",
26             "description" : "spring 在java领域非常流行, java程序员都在用。",
27             "studymodel" : "201001",
28             "price" : 88.6,
29             "timestamp" : "2019-08-24 19:11:35",
30             "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
31             "tags" : [
32                 "spring",
33                 "java"
34             ]
35         },
36         "_explanation" : {
37             "value" : 2.137549,
38             "description" : "sum of:",
39             "details" : [
40                 {
41                     "value" : 0.7936629,
42                     "description" : "weight(description:java in 0)
[PerFieldSimilarity], result of:",
43                     "details" : [
44                         {
45                             "value" : 0.7936629,
46                             "description" : "score(freq=2.0), product of:",
47                             "details" : [
48                                 {
49                                     "value" : 2.2,
50                                     "description" : "boost",
51                                     "details" : [ ]
52                                 },
53                                 {
54                                     "value" : 0.47000363,
55                                     "description" : "idf, computed as log(1 + (N - n + 0.5) /
(n + 0.5)) from:",
56                                     "details" : [
57                                         {
58                                             "value" : 2,
59                                             "description" : "n, number of documents containing
term",
60                                             "details" : [ ]
61                                         },
62                                         {
63                                             "value" : 3,
64                                             "description" : "N, total number of documents with
field",
65                                             "details" : [ ]
66                                         }
67                                     ]
68                                 },
69                                 {
70                                     "value" : 0.7675597,
71                                     "description" : "tf, computed as freq / (freq + k1 * (1 -
b + b * dl / avgdl)) from:",
72                                     "details" : [
73                                         {
74                                             "value" : 2.0,
75                                             "description" : "freq, occurrences of term within
document",

```

```

76         "details" : [ ]
77     },
78     {
79         "value" : 1.2,
80         "description" : "k1, term saturation parameter",
81         "details" : [ ]
82     },
83     {
84         "value" : 0.75,
85         "description" : "b, length normalization parameter",
86         "details" : [ ]
87     },
88     {
89         "value" : 12.0,
90         "description" : "dl, length of field",
91         "details" : [ ]
92     },
93     {
94         "value" : 35.333332,
95         "description" : "avgdl, average length of field",
96         "details" : [ ]
97     }
98 ]
99 }
100 ]
101 }
102 ]
103 },
104 {
105     "value" : 1.3438859,
106     "description" : "weight(description:程序员 in 0)
[PerFieldSimilarity], result of:",
107     "details" : [
108         {
109             "value" : 1.3438859,
110             "description" : "score(freq=1.0), product of:",
111             "details" : [
112                 {
113                     "value" : 2.2,
114                     "description" : "boost",
115                     "details" : [ ]
116                 },
117                 {
118                     "value" : 0.98082924,
119                     "description" : "idf, computed as log(1 + (N - n + 0.5) /
(n + 0.5)) from:",
120                     "details" : [
121                         {
122                             "value" : 1,
123                             "description" : "n, number of documents containing
term",
124                             "details" : [ ]
125                         },
126                         {
127                             "value" : 3,
128                             "description" : "N, total number of documents with
field",
129                             "details" : [ ]

```

```

130         }
131     ]
132 },
133 {
134     "value" : 0.6227967,
135     "description" : "tf, computed as freq / (freq + k1 * (1 -
b + b * dl / avgdl)) from:",
136     "details" : [
137         {
138             "value" : 1.0,
139             "description" : "freq, occurrences of term within
document",
140             "details" : [ ]
141         },
142         {
143             "value" : 1.2,
144             "description" : "k1, term saturation parameter",
145             "details" : [ ]
146         },
147         {
148             "value" : 0.75,
149             "description" : "b, length normalization parameter",
150             "details" : [ ]
151         },
152         {
153             "value" : 12.0,
154             "description" : "dl, length of field",
155             "details" : [ ]
156         },
157         {
158             "value" : 35.333332,
159             "description" : "avgdl, average length of field",
160             "details" : [ ]
161         }
162     ]
163 }
164 ]
165 }
166 ]
167 }
168 ]
169 }
170 },
171 {
172     "_shard" : "[book][0]",
173     "_node" : "MDA45-r6SUGJ0ZyqyhTINA",
174     "_index" : "book",
175     "_type" : "_doc",
176     "_id" : "2",
177     "_score" : 0.57961315,
178     "_source" : {
179         "name" : "java编程思想",
180         "description" : "java语言是世界第一编程语言，在软件开发领域使用人数最多。",
181         "studymodel" : "201001",
182         "price" : 68.6,
183         "timestamp" : "2019-08-25 19:11:35",
184         "pic" : "group1/M00/00/00/wKh1QFs6RCeAY0pHAAJx5ZjNDEM428.jpg",
185         "tags" : [

```

```

186         "java",
187         "dev"
188     ]
189 },
190     "_explanation" : {
191         "value" : 0.57961315,
192         "description" : "sum of:",
193         "details" : [
194             {
195                 "value" : 0.57961315,
196                 "description" : "weight(description:java in 0)
[PerFieldSimilarity], result of:",
197                 "details" : [
198                     {
199                         "value" : 0.57961315,
200                         "description" : "score(freq=1.0), product of:",
201                         "details" : [
202                             {
203                                 "value" : 2.2,
204                                 "description" : "boost",
205                                 "details" : [ ]
206                             },
207                             {
208                                 "value" : 0.47000363,
209                                 "description" : "idf, computed as log(1 + (N - n + 0.5) /
(n + 0.5)) from:",
210                                 "details" : [
211                                     {
212                                         "value" : 2,
213                                         "description" : "n, number of documents containing
term",
214                                         "details" : [ ]
215                                     },
216                                     {
217                                         "value" : 3,
218                                         "description" : "N, total number of documents with
field",
219                                         "details" : [ ]
220                                     }
221                                 ]
222                             },
223                             {
224                                 "value" : 0.56055,
225                                 "description" : "tf, computed as freq / (freq + k1 * (1 -
b + b * dl / avgdl)) from:",
226                                 "details" : [
227                                     {
228                                         "value" : 1.0,
229                                         "description" : "freq, occurrences of term within
document",
230                                         "details" : [ ]
231                                     },
232                                     {
233                                         "value" : 1.2,
234                                         "description" : "k1, term saturation parameter",
235                                         "details" : [ ]
236                                     },
237                                     {

```

```

238         "value" : 0.75,
239         "description" : "b, length normalization parameter",
240         "details" : [ ]
241     },
242     {
243         "value" : 19.0,
244         "description" : "dl, length of field",
245         "details" : [ ]
246     },
247     {
248         "value" : 35.333332,
249         "description" : "avgdl, average length of field",
250         "details" : [ ]
251     }
252 ]
253 }
254 ]
255 }
256 ]
257 }
258 ]
259 }
260 }
261 ]
262 }
263 }

```

16.1.3 分析一个document是如何被匹配上的

```

1  GET /book/_explain/3
2  {
3      "query": {
4          "match": {
5              "description": "java程序员"
6          }
7      }
8  }

```

16.2. Doc value

搜索的时候，要依靠倒排索引；排序的时候，需要依靠正排索引，看到每个document的每个field，然后进行排序，所谓的正排索引，其实就是doc values

在建立索引的时候，一方面会建立倒排索引，以供搜索用；一方面会建立正排索引，也就是doc values，以供排序，聚合，过滤等操作使用

doc values是被保存在磁盘上的，此时如果内存足够，os会自动将其缓存在内存中，性能还是会很高；如果内存不足够，os会将其写入磁盘上

倒排索引

doc1: hello world you and me

doc2: hi, world, how are you

term	doc1	doc2
hello	*	
world	*	*
you	*	*
and	*	
me	*	
hi		*
how		*
are		*

搜索时：

hello you --> hello, you

hello --> doc1

you --> doc1,doc2

doc1: hello world you and me

doc2: hi, world, how are you

sort by 出现问题

正排索引

doc1: { "name": "jack", "age": 27 }

doc2: { "name": "tom", "age": 30 }

document	name	age
doc1	jack	27
doc2	tom	30

#16.3. query phase

#1、query phase

(1) 搜索请求发送到某一个coordinate node，构构建一个priority queue，长度以paging操作from和size为准，默认为10

(2) coordinate node将请求转发到所有shard，每个shard本地搜索，并构建一个本地的priority queue

(3) 各个shard将自己的priority queue返回给coordinate node，并构建一个全局的priority queue

#2、replica shard如何提升搜索吞吐量

一次请求要打到所有shard的一个replica/primary上去，如果每个shard都有多个replica，那么同时并发过来的搜索请求可以同时打到其他的replica上去

#16.4. fetch phase

#1、fetch phase工作流程

- (1) coordinate node构建完priority queue之后，就发送mget请求去所有shard上获取对应的document
- (2) 各个shard将document返回给coordinate node
- (3) coordinate node将合并后的document结果返回给client客户端

#2、一般搜索，如果不加from和size，就默认搜索前10条，按照_score排序

#16.5. 搜索参数小总结

#1、preference

决定了哪些shard会被用来执行搜索操作

_primary, _primary_first, _local, _only_node:xyz, _prefer_node:xyz, _shards:2,3

bouncing results问题，两个document排序，field值相同；不同的shard上，可能排序不同；每次请求轮询打到不同的replica shard上；每次页面上看到的搜索结果的排序都不一样。这就是bouncing result，也就是跳跃的结果。

搜索的时候，是轮询将搜索请求发送到每一个replica shard（primary shard），但是在不同的shard上，可能document的排序不同

解决方案就是将preference设置为一个字符串，比如说user_id，让每个user每次搜索的时候，都使用同一个replica shard去执行，就不会看到bouncing results了

#2、timeout

已经讲解过原理了，主要就是限定在一定时间内，将部分获取到的数据直接返回，避免查询耗时过长

#3、routing

document文档路由，_id路由，routing=user_id，这样的话可以让同一个user对应的数据到一个shard上去

#4、search_type

default: query_then_fetch

dfs_query_then_fetch，可以提升relevance sort精准度

#17. 聚合入门

#17.1聚合示例

#17.1.1需求：计算每个studymodel下的商品数量

sql语句：select studymodel, count(*) from book group by studymodel

```
1  GET /book/_search
2  {
3    "size": 0,
4    "query": {
5      "match_all": {}
6    },
7    "aggs": {
8      "group_by_model": {
9        "terms": { "field": "studymodel" }
10     }
11   }
12 }
```

#17.1.2 需求：计算每个tags下的商品数量

设置字段"fielddata": true

```
1  PUT /book/_mapping/
2  {
3    "properties": {
4      "tags": {
5        "type": "text",
6        "fielddata": true
7      }
8    }
9  }
```

查询

```
1  GET /book/_search
2  {
3    "size": 0,
4    "query": {
5      "match_all": {}
6    },
7    "aggs": {
8      "group_by_tags": {
9        "terms": { "field": "tags" }
10     }
11   }
12 }
```

#17.1.3 需求：加上搜索条件，计算每个tags下的商品数量

```
1  GET /book/_search
2  {
3    "size": 0,
4    "query": {
5      "match": {
```

```

6      "description": "java程序员"
7    }
8  },
9  "aggs": {
10    "group_by_tags": {
11      "terms": { "field": "tags" }
12    }
13  }
14 }

```

#17.1.4 需求：先分组，再算每组的平均值，计算每个tag下的商品的平均价格

```

1  GET /book/_search
2  {
3    "size": 0,
4    "aggs" : {
5      "group_by_tags" : {
6        "terms" : {
7          "field" : "tags"
8        },
9        "aggs" : {
10         "avg_price" : {
11           "avg" : { "field" : "price" }
12         }
13       }
14     }
15   }
16 }

```

#17.1.5 需求：计算每个tag下的商品的平均价格，并且按照平均价格降序排序

```

1  GET /book/_search
2  {
3    "size": 0,
4    "aggs" : {
5      "group_by_tags" : {
6        "terms" : {
7          "field" : "tags",
8          "order": {
9            "avg_price": "desc"
10          }
11        },
12        "aggs" : {
13          "avg_price" : {
14            "avg" : { "field" : "price" }
15          }
16        }
17      }
18    }
19 }

```

#17.1.6 需求：按照指定的价格范围区间进行分组，然后在每组内再按照tag进行分组，最后再计算每组的平均价格

```
1  GET /book/_search
2  {
3    "size": 0,
4    "aggs": {
5      "group_by_price": {
6        "range": {
7          "field": "price",
8          "ranges": [
9            {
10             "from": 0,
11             "to": 40
12           },
13           {
14             "from": 40,
15             "to": 60
16           },
17           {
18             "from": 60,
19             "to": 80
20           }
21         ]
22       },
23       "aggs": {
24         "group_by_tags": {
25           "terms": {
26             "field": "tags"
27           },
28           "aggs": {
29             "average_price": {
30               "avg": {
31                 "field": "price"
32               }
33             }
34           }
35         }
36       }
37     }
38   }
39 }
```

#17.2两个核心概念： bucket和metric

#17.2.1 bucket： 一个数据分组

city name 北京 张三 北京 李四 天津 王五 天津 赵六

天津 王麻子

划分出来两个bucket，一个是北京bucket，一个是天津bucket 北京bucket： 包含了2个人， 张三， 李四 上海bucket： 包含了3个人， 王五， 赵六， 王麻子

#17.2.2metric： 对一个数据分组执行的统计

metric，就是对一个bucket执行的某种聚合分析的操作，比如说求平均值，求最大值，求最小值

```
select count(*) from book group studymodel
```

bucket：group by studymodel --> 那些studymodel相同的数据，就会被划分到一个bucket中 metric：count(*)，对每个user_id bucket中所有的数据，计算一个数量。还有avg(), sum(), max(), min()

#17.3 电视案例

创建索引及映射

```
1  PUT /tvs
2  PUT /tvs/_search
3  {
4      "properties": {
5          "price": {
6              "type": "long"
7          },
8          "color": {
9              "type": "keyword"
10         },
11         "brand": {
12             "type": "keyword"
13         },
14         "sold_date": {
15             "type": "date"
16         }
17     }
18 }
```

插入数据

```
1  POST /tvs/_bulk
2  { "index": {} }
3  { "price" : 1000, "color" : "红色", "brand" : "长虹", "sold_date" : "2019-10-28" }
4  { "index": {} }
5  { "price" : 2000, "color" : "红色", "brand" : "长虹", "sold_date" : "2019-11-05" }
6  { "index": {} }
7  { "price" : 3000, "color" : "绿色", "brand" : "小米", "sold_date" : "2019-05-18" }
8  { "index": {} }
9  { "price" : 1500, "color" : "蓝色", "brand" : "TCL", "sold_date" : "2019-07-02" }
10 { "index": {} }
11 { "price" : 1200, "color" : "绿色", "brand" : "TCL", "sold_date" : "2019-08-19" }
12 { "index": {} }
13 { "price" : 2000, "color" : "红色", "brand" : "长虹", "sold_date" : "2019-11-05" }
14 { "index": {} }
15 { "price" : 8000, "color" : "红色", "brand" : "三星", "sold_date" : "2020-01-01" }
16 { "index": {} }
17 { "price" : 2500, "color" : "蓝色", "brand" : "小米", "sold_date" : "2020-02-12" }
```

#需求1 统计哪种颜色的电视销量最高

```
1  GET /tvs/_search
2  {
3      "size" : 0,
4      "aggs" : {
5          "popular_colors" : {
6              "terms" : {
7                  "field" : "color"
8              }
9          }
10     }
11 }
```

查询条件解析

size: 只获取聚合结果，而不要执行聚合的原始数据 aggs: 固定语法，要对一份数据执行分组聚合操作
popular_colors: 就是对每个aggs，都要起一个名字， terms: 根据字段的值进行分组 field: 根据指定的字段的值进行分组

返回

```
1  {
2      "took" : 18,
3      "timed_out" : false,
4      "_shards" : {
5          "total" : 1,
6          "successful" : 1,
7          "skipped" : 0,
8          "failed" : 0
9      },
10     "hits" : {
11         "total" : {
12             "value" : 8,
13             "relation" : "eq"
14         },
15         "max_score" : null,
16         "hits" : [ ]
17     },
18     "aggregations" : {
19         "popular_colors" : {
20             "doc_count_error_upper_bound" : 0,
21             "sum_other_doc_count" : 0,
22             "buckets" : [
23                 {
24                     "key" : "红色",
25                     "doc_count" : 4
26                 },
27                 {
28                     "key" : "绿色",
29                     "doc_count" : 2
30                 },
31                 {
32                     "key" : "蓝色",
33                     "doc_count" : 2
34                 }
35             ]
36         }
37     }
```

```
37     }
38 }
```

返回结果解析

hits.hits: 我们指定了size是0, 所以hits.hits就是空的 aggregations: 聚合结果 popular_color: 我们指定的某个聚合的名称 buckets: 根据我们指定的field划分出的buckets key: 每个bucket对应的那个值

doc_count: 这个bucket分组内, 有多少个数据 数量, 其实就是这种颜色的销量

每种颜色对应的bucket中的数据的默认的排序规则: 按照doc_count降序排序

#需求2 统计每种颜色电视平均价格

```
1  GET /tvs/_search
2  {
3    "size" : 0,
4    "aggs": {
5      "colors": {
6        "terms": {
7          "field": "color"
8        },
9        "aggs": {
10       "avg_price": {
11         "avg": {
12           "field": "price"
13         }
14       }
15     }
16   }
17 }
18 }
```

在一个aggs执行的bucket操作 (terms), 平级的json结构下, 再加一个aggs, 这个第二个aggs内部, 同样取个名字, 执行一个metric操作, avg, 对之前的每个bucket中的数据的指定的field, price field, 求一个平均值

返回:

```
1  {
2    "took" : 4,
3    "timed_out" : false,
4    "_shards" : {
5      "total" : 1,
6      "successful" : 1,
7      "skipped" : 0,
8      "failed" : 0
9    },
10   "hits" : {
11     "total" : {
12       "value" : 8,
13       "relation" : "eq"
14     },
15     "max_score" : null,
16     "hits" : [ ]
17   },
18   "aggregations" : {
19     "colors" : {
20       "doc_count_error_upper_bound" : 0,
```



```

21     "sum_other_doc_count" : 0,
22     "buckets" : [
23       {
24         "key" : "红色",
25         "doc_count" : 4,
26         "avg_price" : {
27           "value" : 3250.0
28         }
29       },
30       {
31         "key" : "绿色",
32         "doc_count" : 2,
33         "avg_price" : {
34           "value" : 2100.0
35         }
36       },
37       {
38         "key" : "蓝色",
39         "doc_count" : 2,
40         "avg_price" : {
41           "value" : 2000.0
42         }
43       }
44     ]
45   }
46 }
47 }

```

buckets, 除了key和doc_count avg_price: 我们自己取的metric aggs的名字 value: 我们的metric计算的结果, 每个bucket中的数据的价格字段求平均值后的结果

相当于sql: select avg(price) from tvs group by color

#需求3 继续下钻分析

每个颜色下, 平均价格及每个颜色下, 每个品牌的平均价格

```

1  GET /tvs/_search
2  {
3    "size": 0,
4    "aggs": {
5      "group_by_color": {
6        "terms": {
7          "field": "color"
8        },
9        "aggs": {
10         "color_avg_price": {
11           "avg": {
12             "field": "price"
13           }
14         },
15         "group_by_brand": {
16           "terms": {
17             "field": "brand"
18           },
19           "aggs": {
20             "brand_avg_price": {
21               "avg": {

```

```

22         "field": "price"
23     }
24 }
25 }
26 }
27 }
28 }
29 }
30 }

```

#需求4：更多的metric

count: bucket, terms, 自动就会有一个doc_count, 就相当于count avg: avg aggs, 求平均值 max: 求一个bucket内, 指定field值最大的那个数据 min: 求一个bucket内, 指定field值最小的那个数据 sum: 求一个bucket内, 指定field值的总和

```

1  GET /tvs/_search
2  {
3      "size" : 0,
4      "aggs": {
5          "colors": {
6              "terms": {
7                  "field": "color"
8              },
9              "aggs": {
10                 "avg_price": { "avg": { "field": "price" } },
11                 "min_price" : { "min": { "field": "price" } },
12                 "max_price" : { "max": { "field": "price" } },
13                 "sum_price" : { "sum": { "field": "price" } }
14             }
15         }
16     }
17 }

```

#需求5：划分范围 histogram

```

1  GET /tvs/_search
2  {
3      "size" : 0,
4      "aggs":{
5          "price":{
6              "histogram":{
7                  "field": "price",
8                  "interval": 2000
9              },
10             "aggs":{
11                 "income": {
12                     "sum": {
13                         "field" : "price"
14                     }
15                 }
16             }
17         }
18     }
19 }

```

histogram：类似于terms，也是进行bucket分组操作，接收一个field，按照这个field的值的各个范围区间，进行bucket分组操作

```
1  "histogram":{
2    "field": "price",
3    "interval": 2000
4  }
```

interval: 2000, 划分范围, 0_2000, 2000_4000, 4000_6000, 6000_8000, 8000~10000, buckets

bucket有了之后，一样的，去对每个bucket执行avg, count, sum, max, min, 等各种metric操作，聚合分析

#需求6：按照日期分组聚合

date_histogram，按照我们指定的某个date类型的日期field，以及日期interval，按照一定的日期间隔，去划分bucket

min_doc_count：即使某个日期interval，2017-01-01~2017-01-31中，一条数据都没有，那么这个区间也是要返回的，不然默认是会过滤掉这个区间的 extended_bounds, min, max：划分bucket的时候，会限定在这个起始日期，和截止日期内

```
1  GET /tvs/_search
2  {
3    "size" : 0,
4    "aggs": {
5      "sales": {
6        "date_histogram": {
7          "field": "sold_date",
8          "interval": "month",
9          "format": "yyyy-MM-dd",
10         "min_doc_count" : 0,
11         "extended_bounds" : {
12           "min" : "2019-01-01",
13           "max" : "2020-12-31"
14         }
15       }
16     }
17   }
18 }
```

#需求7 统计每季度每个品牌的销售额

```
1  GET /tvs/_search
2  {
3    "size": 0,
4    "aggs": {
5      "group_by_sold_date": {
6        "date_histogram": {
7          "field": "sold_date",
8          "interval": "quarter",
9          "format": "yyyy-MM-dd",
10         "min_doc_count": 0,
11         "extended_bounds": {
12           "min": "2019-01-01",
13           "max": "2020-12-31"
14         }
15       }
16     }
17   }
18 }
```

```

15     },
16     "aggs": {
17         "group_by_brand": {
18             "terms": {
19                 "field": "brand"
20             },
21             "aggs": {
22                 "sum_price": {
23                     "sum": {
24                         "field": "price"
25                     }
26                 }
27             }
28         },
29         "total_sum_price": {
30             "sum": {
31                 "field": "price"
32             }
33         }
34     }
35 }
36 }
37 }

```

求8：搜索与聚合结合，查询某个品牌按颜色销量

搜索与聚合可以结合起来。

sql select count(*)

from tvs

where brand like "%小米%"

group by color

es aggregation, scope, 任何的聚合，都必须在搜索出来的结果数据中之行，搜索结果，就是聚合分析操作的scope

```

1  GET /tvs/_search
2  {
3    "size": 0,
4    "query": {
5      "term": {
6        "brand": {
7          "value": "小米"
8        }
9      }
10   },
11   "aggs": {
12     "group_by_color": {
13       "terms": {
14         "field": "color"
15       }
16     }
17   }
18 }

```

#需求9 global bucket：单个品牌与所有品牌销量对比

aggregation, scope, 一个聚合操作, 必须在query的搜索结果范围内执行

出来两个结果, 一个结果, 是基于query搜索结果来聚合的; 一个结果, 是对所有数据执行聚合的

```
1  GET /tvs/_search
2  {
3    "size": 0,
4    "query": {
5      "term": {
6        "brand": {
7          "value": "小米"
8        }
9      }
10   },
11   "aggs": {
12     "single_brand_avg_price": {
13       "avg": {
14         "field": "price"
15       }
16     },
17     "all": {
18       "global": {},
19       "aggs": {
20         "all_brand_avg_price": {
21           "avg": {
22             "field": "price"
23           }
24         }
25       }
26     }
27   }
28 }
```

#需求10: 过滤+聚合：统计价格大于1200的电视平均价格

搜索+聚合

过滤+聚合

```
1  GET /tvs/_search
2  {
3    "size": 0,
4    "query": {
5      "constant_score": {
6        "filter": {
7          "range": {
8            "price": {
9              "gte": 1200
10           }
11         }
12       }
13     },
14   },
15   "aggs": {
16     "avg_price": {
17       "avg": {
```

```
18         "field": "price"
19     }
20 }
21 }
22 }
```

#需求11 bucket filter: 统计品牌最近一个月的平均价格

```
1  GET /tvs/_search
2  {
3      "size": 0,
4      "query": {
5          "term": {
6              "brand": {
7                  "value": "小米"
8              }
9          }
10 },
11 "aggs": {
12     "recent_150d": {
13         "filter": {
14             "range": {
15                 "sold_date": {
16                     "gte": "now-150d"
17                 }
18             }
19         },
20         "aggs": {
21             "recent_150d_avg_price": {
22                 "avg": {
23                     "field": "price"
24                 }
25             }
26         }
27     },
28     "recent_140d": {
29         "filter": {
30             "range": {
31                 "sold_date": {
32                     "gte": "now-140d"
33                 }
34             }
35         },
36         "aggs": {
37             "recent_140d_avg_price": {
38                 "avg": {
39                     "field": "price"
40                 }
41             }
42         }
43     },
44     "recent_130d": {
45         "filter": {
46             "range": {
47                 "sold_date": {
48                     "gte": "now-130d"
49                 }
50             }
51         },
52         "aggs": {
53             "recent_130d_avg_price": {
54                 "avg": {
55                     "field": "price"
56                 }
57             }
58         }
59     }
60 }
```

```

50     }
51   },
52   "aggs": {
53     "recent_130d_avg_price": {
54       "avg": {
55         "field": "price"
56       }
57     }
58   }
59 }
60 }
61 }

```

aggs.filter, 针对的是聚合去做的

如果放query里面的filter, 是全局的, 会对所有的数据都有影响

但是, 如果, 比如说, 你要统计, 长虹电视, 最近1个月的平均值; 最近3个月的平均值; 最近6个月的平均值

bucket filter: 对不同的bucket下的aggs, 进行filter

需求12 排序: 按每种颜色的平均销售额降序排序

```

1  GET /tvs/_search
2  {
3    "size": 0,
4    "aggs": {
5      "group_by_color": {
6        "terms": {
7          "field": "color",
8          "order": {
9            "avg_price": "asc"
10         }
11       },
12       "aggs": {
13         "avg_price": {
14           "avg": {
15             "field": "price"
16           }
17         }
18       }
19     }
20   }
21 }

```

相当于sql子表数据字段可以立刻使用。

需求13 排序: 按每种颜色的每种品牌平均销售额降序排序

```

1  GET /tvs/_search
2  {
3    "size": 0,
4    "aggs": {
5      "group_by_color": {
6        "terms": {
7          "field": "color"
8        },

```

```

9      "aggs": {
10        "group_by_brand": {
11          "terms": {
12            "field": "brand",
13            "order": {
14              "avg_price": "desc"
15            }
16          },
17          "aggs": {
18            "avg_price": {
19              "avg": {
20                "field": "price"
21              }
22            }
23          }
24        }
25      }
26    }
27  }
28 }

```

#18. java api实现聚合

简单聚合，多种聚合，详见代码。

```

1  package com.ydlclass.es;
2
3  import org.elasticsearch.action.search.SearchRequest;
4  import org.elasticsearch.action.search.SearchResponse;
5  import org.elasticsearch.client.RequestOptions;
6  import org.elasticsearch.client.RestHighLevelClient;
7  import org.elasticsearch.index.query.QueryBuilders;
8  import org.elasticsearch.search.aggregations.Aggregation;
9  import org.elasticsearch.search.aggregations.AggregationBuilders;
10 import org.elasticsearch.search.aggregations.Aggregations;
11 import org.elasticsearch.search.aggregations.bucket.histogram.*;
12 import org.elasticsearch.search.aggregations.bucket.terms.Terms;
13 import
    org.elasticsearch.search.aggregations.bucket.terms.TermsAggregationBuilder;
14 import org.elasticsearch.search.aggregations.metrics.*;
15 import org.elasticsearch.search.builder.SearchSourceBuilder;
16 import org.junit.Test;
17 import org.junit.runner.RunWith;
18 import org.springframework.beans.factory.annotation.Autowired;
19 import org.springframework.boot.test.context.SpringBootTest;
20 import org.springframework.test.context.junit4.SpringRunner;
21
22 import java.io.IOException;
23 import java.util.List;
24
25 /**
26  * creste by ydlclass.ydl
27  */
28 @SpringBootTest
29 @RunWith(SpringRunner.class)
30 public class TestAggs {
31     @Autowired

```



```

32     RestHighLevelClient client;
33
34     //需求一：按照颜色分组，计算每个颜色卖出的个数
35     @Test
36     public void testAggs() throws IOException {
37         // GET /tvs/_search
38         // {
39         //     "size": 0,
40         //     "query": {"match_all": {}},
41         //     "aggs": {
42         //         "group_by_color": {
43         //             "terms": {
44         //                 "field": "color"
45         //             }
46         //         }
47         //     }
48         // }
49
50         //1 构建请求
51         SearchRequest searchRequest=new SearchRequest("tvs");
52
53         //请求体
54         SearchSourceBuilder searchSourceBuilder=new SearchSourceBuilder();
55         searchSourceBuilder.size(0);
56         searchSourceBuilder.query(QueryBuilders.matchAllQuery());
57
58         TermsAggregationBuilder termsAggregationBuilder =
59 AggregationBuilders.terms("group_by_color").field("color");
60         searchSourceBuilder.aggregation(termsAggregationBuilder);
61
62         //请求体放入请求头
63         searchRequest.source(searchSourceBuilder);
64
65         //2 执行
66         SearchResponse searchResponse = client.search(searchRequest,
67 RequestOptions.DEFAULT);
68
69         //3 获取结果
70         // "aggregations" : {
71         //     "group_by_color" : {
72         //         "doc_count_error_upper_bound" : 0,
73         //         "sum_other_doc_count" : 0,
74         //         "buckets" : [
75         //             {
76         //                 "key" : "红色",
77         //                 "doc_count" : 4
78         //             },
79         //             {
80         //                 "key" : "绿色",
81         //                 "doc_count" : 2
82         //             },
83         //             {
84         //                 "key" : "蓝色",
85         //                 "doc_count" : 2
86         //             }
87         //         ]
88         //     }
89         // }
90         Aggregations aggregations = searchResponse.getAggregations();

```

```

88     Terms group_by_color = aggregations.get("group_by_color");
89     List<? extends Terms.Bucket> buckets = group_by_color.getBuckets();
90     for (Terms.Bucket bucket : buckets) {
91         String key = bucket.getKeyAsString();
92         System.out.println("key:"+key);
93
94         long docCount = bucket.getDocCount();
95         System.out.println("docCount:"+docCount);
96
97         System.out.println("=====");
98     }
99 }
100
101 // #需求二：按照颜色分组，计算每个颜色卖出的个数，每个颜色卖出的平均价格
102 @Test
103 public void testAggsAndAvg() throws IOException {
104     // GET /tvs/_search
105     // {
106     //     "size": 0,
107     //     "query": {"match_all": {}},
108     //     "aggs": {
109     //         "group_by_color": {
110     //             "terms": {
111     //                 "field": "color"
112     //             },
113     //             "aggs": {
114     //                 "avg_price": {
115     //                     "avg": {
116     //                         "field": "price"
117     //                     }
118     //                 }
119     //             }
120     //         }
121     //     }
122     // }
123
124     //1 构建请求
125     SearchRequest searchRequest=new SearchRequest("tvs");
126
127     //请求体
128     SearchSourceBuilder searchSourceBuilder=new SearchSourceBuilder();
129     searchSourceBuilder.size(0);
130     searchSourceBuilder.query(QueryBuilders.matchAllQuery());
131
132     TermsAggregationBuilder termsAggregationBuilder =
133     AggregationBuilders.terms("group_by_color").field("color");
134
135     //terms聚合下填充一个子聚合
136     AvgAggregationBuilder avgAggregationBuilder =
137     AggregationBuilders.avg("avg_price").field("price");
138     termsAggregationBuilder.subAggregation(avgAggregationBuilder);
139
140     searchSourceBuilder.aggregation(termsAggregationBuilder);
141
142     //请求体放入请求头
143     searchRequest.source(searchSourceBuilder);

```

```

144         SearchResponse searchResponse = client.search(searchRequest,
RequestOptions.DEFAULT);

145
146         //3 获取结果
147         // {
148         //     "key" : "红色",
149         //     "doc_count" : 4,
150         //     "avg_price" : {
151         //         "value" : 3250.0
152         //     }
153         // }
154         Aggregations aggregations = searchResponse.getAggregations();
155         Terms group_by_color = aggregations.get("group_by_color");
156         List<? extends Terms.Bucket> buckets = group_by_color.getBuckets();
157         for (Terms.Bucket bucket : buckets) {
158             String key = bucket.getKeyAsString();
159             System.out.println("key:"+key);
160
161             long docCount = bucket.getDocCount();
162             System.out.println("docCount:"+docCount);
163
164             Aggregations aggregations1 = bucket.getAggregations();
165             Avg avg_price = aggregations1.get("avg_price");
166             double value = avg_price.getValue();
167             System.out.println("value:"+value);
168
169             System.out.println("=====");
170         }
171     }
172
173     // #需求三: 按照颜色分组, 计算每个颜色卖出的个数, 以及每个颜色卖出的平均值、最大值、最小
值、总和。
174     @Test
175     public void testAggsAndMore() throws IOException {
176         // GET /tvs/_search
177         // {
178         //     "size" : 0,
179         //     "aggs": {
180         //         "group_by_color": {
181         //             "terms": {
182         //                 "field": "color"
183         //             },
184         //             "aggs": {
185         //                 "avg_price": { "avg": { "field": "price" } },
186         //                 "min_price" : { "min": { "field": "price" } },
187         //                 "max_price" : { "max": { "field": "price" } },
188         //                 "sum_price" : { "sum": { "field": "price" } }
189         //             }
190         //         }
191         //     }
192         // }
193
194         //1 构建请求
195         SearchRequest searchRequest=new SearchRequest("tvs");
196
197         //请求体
198         SearchSourceBuilder searchSourceBuilder=new SearchSourceBuilder();
199         searchSourceBuilder.size(0);

```

```

200         searchSourceBuilder.query(QueryBuilders.matchAllQuery());
201
202         TermsAggregationBuilder termsAggregationBuilder =
AggregationBuilders.terms("group_by_color").field("color");
203
204
205         //termsAggregationBuilder里放入多个子聚合
206         AvgAggregationBuilder avgAggregationBuilder =
AggregationBuilders.avg("avg_price").field("price");
207         MinAggregationBuilder minAggregationBuilder =
AggregationBuilders.min("min_price").field("price");
208         MaxAggregationBuilder maxAggregationBuilder =
AggregationBuilders.max("max_price").field("price");
209         SumAggregationBuilder sumAggregationBuilder =
AggregationBuilders.sum("sum_price").field("price");
210
211         termsAggregationBuilder.subAggregation(avgAggregationBuilder);
212         termsAggregationBuilder.subAggregation(minAggregationBuilder);
213         termsAggregationBuilder.subAggregation(maxAggregationBuilder);
214         termsAggregationBuilder.subAggregation(sumAggregationBuilder);
215
216
217         searchSourceBuilder.aggregation(termsAggregationBuilder);
218
219         //请求体放入请求头
220         searchRequest.source(searchSourceBuilder);
221
222         //2 执行
223         SearchResponse searchResponse = client.search(searchRequest,
RequestOptions.DEFAULT);
224
225         //3 获取结果
226         // {
227         //     "key" : "红色",
228         //     "doc_count" : 4,
229         //     "max_price" : {
230         //         "value" : 8000.0
231         //     },
232         //     "min_price" : {
233         //         "value" : 1000.0
234         //     },
235         //     "avg_price" : {
236         //         "value" : 3250.0
237         //     },
238         //     "sum_price" : {
239         //         "value" : 13000.0
240         //     }
241         // }
242         Aggregations aggregations = searchResponse.getAggregations();
243         Terms group_by_color = aggregations.get("group_by_color");
244         List<? extends Terms.Bucket> buckets = group_by_color.getBuckets();
245         for (Terms.Bucket bucket : buckets) {
246             String key = bucket.getKeyAsString();
247             System.out.println("key:"+key);
248
249             long docCount = bucket.getDocCount();
250             System.out.println("docCount:"+docCount);
251

```

```

252         Aggregations aggregations1 = bucket.getAggregations();
253
254         Max max_price = aggregations1.get("max_price");
255         double maxPriceValue = max_price.getValue();
256         System.out.println("maxPriceValue:"+maxPriceValue);
257
258         Min min_price = aggregations1.get("min_price");
259         double minPriceValue = min_price.getValue();
260         System.out.println("minPriceValue:"+minPriceValue);
261
262         Avg avg_price = aggregations1.get("avg_price");
263         double avgPriceValue = avg_price.getValue();
264         System.out.println("avgPriceValue:"+avgPriceValue);
265
266         Sum sum_price = aggregations1.get("sum_price");
267         double sumPriceValue = sum_price.getValue();
268         System.out.println("sumPriceValue:"+sumPriceValue);
269
270         System.out.println("=====");
271     }
272 }
273
274 // #需求四：按照售价每2000价格划分范围，算出每个区间的销售总额 histogram
275 @Test
276 public void testAggsAndHistogram() throws IOException {
277     // GET /tvs/_search
278     // {
279     //     "size" : 0,
280     //     "aggs":{
281     //         "by_histogram":{
282     //             "histogram":{
283     //                 "field": "price",
284     //                 "interval": 2000
285     //             },
286     //             "aggs":{
287     //                 "income": {
288     //                     "sum": {
289     //                         "field" : "price"
290     //                     }
291     //                 }
292     //             }
293     //         }
294     //     }
295     // }
296
297     //1 构建请求
298     SearchRequest searchRequest=new SearchRequest("tvs");
299
300     //请求体
301     SearchSourceBuilder searchSourceBuilder=new SearchSourceBuilder();
302     searchSourceBuilder.size(0);
303     searchSourceBuilder.query(QueryBuilders.matchAllQuery());
304
305     HistogramAggregationBuilder histogramAggregationBuilder =
306     AggregationBuilders.histogram("by_histogram").field("price").interval(2000);
307
308     SumAggregationBuilder sumAggregationBuilder =
309     AggregationBuilders.sum("income").field("price");

```

```

308     histogramAggregationBuilder.subAggregation(sumAggregationBuilder);
309     searchSourceBuilder.aggregation(histogramAggregationBuilder);
310
311     //请求体放入请求头
312     searchRequest.source(searchSourceBuilder);
313
314     //2 执行
315     SearchResponse searchResponse = client.search(searchRequest,
RequestOptions.DEFAULT);
316
317     //3 获取结果
318     // {
319     //     "key" : 0.0,
320     //     "doc_count" : 3,
321     //     "income" : {
322     //         "value" : 3700.0
323     //     }
324     // }
325     Aggregations aggregations = searchResponse.getAggregations();
326     Histogram group_by_color = aggregations.get("by_histogram");
327     List<? extends Histogram.Bucket> buckets = group_by_color.getBuckets();
328     for (Histogram.Bucket bucket : buckets) {
329         String keyAsString = bucket.getKeyAsString();
330         System.out.println("keyAsString:"+keyAsString);
331         long docCount = bucket.getDocCount();
332         System.out.println("docCount:"+docCount);
333
334         Aggregations aggregations1 = bucket.getAggregations();
335         Sum income = aggregations1.get("income");
336         double value = income.getValue();
337         System.out.println("value:"+value);
338
339         System.out.println("=====");
340
341     }
342 }
343
344 // #需求五: 计算每个季度的销售总额
345 @Test
346 public void testAggsAndDateHistogram() throws IOException {
347     // GET /tvs/_search
348     // {
349     //     "size" : 0,
350     //     "aggs": {
351     //         "sales": {
352     //             "date_histogram": {
353     //                 "field": "sold_date",
354     //                 "interval": "quarter",
355     //                 "format": "yyyy-MM-dd",
356     //                 "min_doc_count" : 0,
357     //                 "extended_bounds" : {
358     //                     "min" : "2019-01-01",
359     //                     "max" : "2020-12-31"
360     //                 }
361     //             },
362     //             "aggs": {
363     //                 "income": {
364     //                     "sum": {

```

```

365         //             "field": "price"
366         //             }
367         //         }
368         //     }
369         // }
370     // }
371 }
372
373 //1 构建请求
374 SearchRequest searchRequest=new SearchRequest("tvs");
375
376 //请求体
377 SearchSourceBuilder searchSourceBuilder=new SearchSourceBuilder();
378 searchSourceBuilder.size(0);
379 searchSourceBuilder.query(QueryBuilders.matchAllQuery());
380
381 DateHistogramAggregationBuilder dateHistogramAggregationBuilder =
AggregationBuilders.dateHistogram("date_histogram").field("sold_date").calendarI
nterval(DateHistogramInterval.QUARTER)
382     .format("yyyy-MM-dd").minDocCount(0).extendedBounds(new
ExtendedBounds("2019-01-01", "2020-12-31"));
383 SumAggregationBuilder sumAggregationBuilder =
AggregationBuilders.sum("income").field("price");
384 dateHistogramAggregationBuilder.subAggregation(sumAggregationBuilder);
385
386 searchSourceBuilder.aggregation(dateHistogramAggregationBuilder);
387 //请求体放入请求头
388 searchRequest.source(searchSourceBuilder);
389
390 //2 执行
391 SearchResponse searchResponse = client.search(searchRequest,
RequestOptions.DEFAULT);
392
393 //3 获取结果
394 // {
395 //     "key_as_string" : "2019-01-01",
396 //     "key" : 1546300800000,
397 //     "doc_count" : 0,
398 //     "income" : {
399 //         "value" : 0.0
400 //     }
401 // }
402 Aggregations aggregations = searchResponse.getAggregations();
403 ParsedDateHistogram date_histogram = aggregations.get("date_histogram");
404 List<? extends Histogram.Bucket> buckets = date_histogram.getBuckets();
405 for (Histogram.Bucket bucket : buckets) {
406     String keyAsString = bucket.getKeyAsString();
407     System.out.println("keyAsString:"+keyAsString);
408     long docCount = bucket.getDocCount();
409     System.out.println("docCount:"+docCount);
410
411     Aggregations aggregations1 = bucket.getAggregations();
412     Sum income = aggregations1.get("income");
413     double value = income.getValue();
414     System.out.println("value:"+value);
415
416     System.out.println("=====");
417 }

```

```
418
419     }
420
421 }
```

#19. es7 sql新特性

#19.1 快速入门

```
1  POST  /_sql?format=txt
2  {
3      "query": "SELECT * FROM tvs "
4  }
```

#19.2 启动方式

1http 请求

2客户端: elasticsearch-sql-cli.bat

3代码

#19.3 显示方式

format	Accept HTTP header	Description
Human Readable		
csv	text/csv	Comma-separated values
json	application/json	JSON (JavaScript Object Notation) human-readable format
tsv	text/tab-separated-values	Tab-separated values
txt	text/plain	CLI-like representation
yaml	application/yaml	YAML (YAML Ain't Markup Language) human-readable format
Binary Formats		
cbor	application/cbor	Concise Binary Object Representation
smile	application/smile	Smile binary data format similar to CBOR

#19.4 sql 翻译

```

1  POST /_sql/translate
2  {
3      "query": "SELECT * FROM tvs "
4  }

```

返回:

```

1  {
2      "size" : 1000,
3      "_source" : false,
4      "stored_fields" : "_none_",
5      "docvalue_fields" : [
6          {
7              "field" : "brand"
8          },
9          {
10             "field" : "color"
11         },
12         {
13             "field" : "price"
14         },
15         {
16             "field" : "sold_date",
17             "format" : "epoch_millis"

```

```

18     }
19   ],
20   "sort" : [
21     {
22       "_doc" : {
23         "order" : "asc"
24       }
25     }
26   ]
27 }

```

#19.5 与其他DSL结合

```

1  POST /_sql?format=txt
2  {
3    "query": "SELECT * FROM tvs",
4    "filter": {
5      "range": {
6        "price": {
7          "gte" : 1200,
8          "lte" : 2000
9        }
10     }
11   }
12 }

```

#19.6 java 代码实现sql功能

1前提 es拥有白金版功能

kibana中管理-》许可管理 开启白金版试用

2导入依赖

```

1      <dependency>
2        <groupId>org.elasticsearch.plugin</groupId>
3        <artifactId>x-pack-sql-jdbc</artifactId>
4        <version>7.3.0</version>
5      </dependency>
6
7      <repositories>
8        <repository>
9          <id>elastic.co</id>
10         <url>https://artifacts.elastic.co/maven</url>
11       </repository>
12     </repositories>

```

3代码

```

1  public static void main(String[] args) {
2    try {
3      Connection connection =
4      DriverManager.getConnection("jdbc:es://http://localhost:9200");
5      Statement statement = connection.createStatement();

```

```

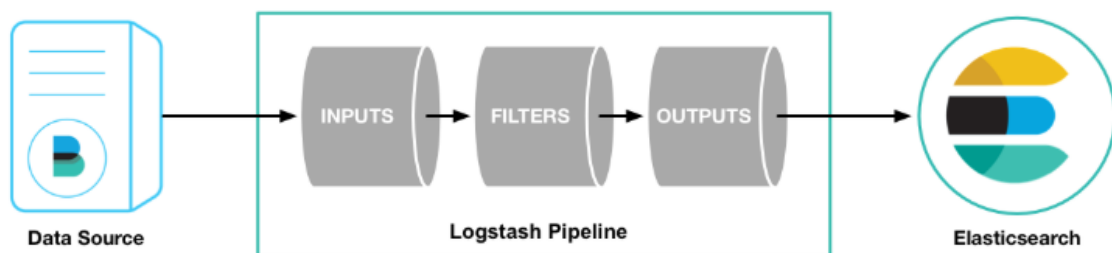
5         ResultSet results = statement.executeQuery(
6             "select * from tvs");
7         while(results.next()){
8             System.out.println(results.getString(1));
9             System.out.println(results.getString(2));
10            System.out.println(results.getString(3));
11            System.out.println(results.getString(4));
12            System.out.println("=====");
13        }
14    }catch (Exception e){
15        e.printStackTrace();
16    }

```

大型企业可以购买白金版，增加Machine Learning、高级安全性x-pack。

#20. Logstash学习

#20.1 Logstash基本语法组成



#1什么是Logstash

logstash是一个数据抽取工具，将数据从一个地方转移到另一个地方。如hadoop生态圈的sqoop等。下载地址：<https://www.elastic.co/cn/downloads/logstash>

logstash之所以功能强大和流行，还与其丰富的过滤器插件是分不开的，过滤器提供的并不单单是过滤的功能，还可以对进入过滤器的原始数据进行复杂的逻辑处理，甚至添加独特的事件到后续流程中。Logstash配置文件有如下三部分组成，其中input、output部分是必须配置，filter部分是可选配置，而filter就是过滤器插件，可以在这部分实现各种日志过滤功能。

#2配置文件：

```

1  input {
2      #输入插件
3  }
4  filter {
5      #过滤匹配插件
6  }
7  output {
8      #输出插件
9  }

```

#3启动操作:

```
1 logstash.bat -e 'input{stdin{}} output{stdout{}}'
```

为了好维护，将配置写入文件，启动

```
1 logstash.bat -f ../config/test1.conf
```

#20.2 Logstash输入插件 (input)

<https://www.elastic.co/guide/en/logstash/current/input-plugins.html>

#1、标准输入(Stdin)

```
1 input{
2     stdin{
3
4     }
5 }
6 output {
7     stdout{
8         codec=>rubydebug
9     }
10 }
```

#2、读取文件(File)

logstash使用一个名为filewatch的ruby gem库来监听文件变化,并通过一个叫.sinceb数据库文件来记录被监听的日志文件的读取进度（时间戳），这个sinceb数据文件的默认路径在<path.data>/plugins/inputs/file下面，文件名类似于.sinceb_123456，而<path.data>表示logstash插件存储目录，默认是LOGSTASH_HOME/data。

```
1 input {
2     file {
3         path => ["/var/*/*"]
4         start_position => "beginning"
5     }
6 }
7 output {
8     stdout{
9         codec=>rubydebug
10    }
11 }
```

默认情况下，logstash会从文件的结束位置开始读取数据，也就是说logstash进程会以类似tail -f命令的形式逐行获取数据。

#3、读取TCP网络数据

```
1 input {
2     tcp {
3         port => "1234"
4     }
5 }
```

```

6
7   filter {
8     grok {
9       match => { "message" => "%{SYSLOGLINE}" }
10    }
11  }
12
13  output {
14    stdout{
15      codec=>rubydebug
16    }
17  }

```

20.3 Logstash过滤器插件(Filter)

<https://www.elastic.co/guide/en/logstash/current/filter-plugins.html>

20.13.1Grok 正则捕获

grok是一个十分强大的logstash filter插件，他可以通过正则解析任意文本，将非结构化日志数据弄成结构化和方便查询的结构。他是目前logstash 中解析非结构化日志数据最好的方式。

Grok 的语法规则是：

```
1   %{语法: 语义}
```

例如输入的内容为：

```
1   172.16.213.132 [07/Feb/2019:16:24:19 +0800] "GET / HTTP/1.1" 403 5039
```

%{IP:clientip}匹配模式将获得的结果为：clientip: 172.16.213.132 %{HTTPDATE:timestamp}匹配模式将获得的结果为：timestamp: 07/Feb/2018:16:24:19 +0800 而%{QS:referrer}匹配模式将获得的结果为：referrer: "GET / HTTP/1.1"

下面是一个组合匹配模式，它可以获取上面输入的所有内容：

```
1   %{IP:clientip}\ \[%{HTTPDATE:timestamp}\]\ %{QS:referrer}\ %{NUMBER:response}\ %
    {NUMBER:bytes}
```

通过上面这个组合匹配模式，我们将输入的内容分成了五个部分，即五个字段，将输入内容分割为不同的数据字段，这对于日后解析和查询日志数据非常有用，这正是使用grok的目的。

例子：

```

1  input{
2      stdin{}
3  }
4  filter{
5      grok{
6          match => ["message", "%{IP:clientip}\ \[%{HTTPDATE:timestamp}\]\ %
          {QS:referrer}\ %{NUMBER:response}\ %{NUMBER:bytes}"]
7      }
8  }
9  output{
10     stdout{
11         codec => "rubydebug"
12     }
13 }

```

输入内容：

```
1 172.16.213.132 [07/Feb/2019:16:24:19 +0800] "GET / HTTP/1.1" 403 5039
```

20.13.2 时间处理(Date)

date插件是对于排序事件和回填旧数据尤其重要，它可以用来转换日志记录中的时间字段，变成LogStash::Timestamp对象，然后转存到@timestamp字段里，这在之前已经做过简单的介绍。下面是date插件的一个配置示例（这里仅仅列出filter部分）：

```

1  filter {
2      grok {
3          match => ["message", "%{HTTPDATE:timestamp}"]
4      }
5      date {
6          match => ["timestamp", "dd/MMM/yyyy:HH:mm:ss Z"]
7      }
8  }

```

20.13.3、数据修改(Mutate)

(1) 正则表达式替换匹配字段

gsub可以通过正则表达式替换字段中匹配到的值，只对字符串字段有效，下面是一个关于mutate插件中gsub的示例（仅列出filter部分）：

```

1  filter {
2      mutate {
3          gsub => ["filed_name_1", "/", "_"]
4      }
5  }

```

这个示例表示将filed_name_1字段中所有"/"字符替换为"_"。

(2) 分隔符分割字符串为数组

split可以通过指定的分隔符分割字段中的字符串为数组，下面是一个关于mutate插件中split的示例（仅列出filter部分）：

```

1   filter {
2       mutate {
3           split => ["filed_name_2", "|"]
4       }
5   }

```

这个示例表示将filed_name_2字段以"|"为区间分隔为数组。

(3) 重命名字段

rename可以实现重命名某个字段的功能，下面是一个关于mutate插件中rename的示例（仅列出filter部分）：

```

1   filter {
2       mutate {
3           rename => { "old_field" => "new_field" }
4       }
5   }

```

这个示例表示将字段old_field重命名为new_field。

(4) 删除字段

remove_field可以实现删除某个字段的功能，下面是一个关于mutate插件中remove_field的示例（仅列出filter部分）：

```

1   filter {
2       mutate {
3           remove_field => ["timestamp"]
4       }
5   }

```

这个示例表示将字段timestamp删除。

(5) GeoIP 地址查询归类

```

1   filter {
2       geoip {
3           source => "ip_field"
4       }
5   }

```

综合例子：

```

1   input {
2       stdin {}
3   }
4   filter {
5       grok {
6           match => { "message" => "%{IP:clientip}\ \[%{HTTPDATE:timestamp}\%\ %
7               {QS:referrer}\ %{NUMBER:response}\ %{NUMBER:bytes}" }
8           remove_field => [ "message" ]
9       }
10      date {
11          match => ["timestamp", "dd/MMM/yyyy:HH:mm:ss Z"]
12      }
13      mutate {
14          convert => [ "response","float" ]
15      }
16  }

```

```

14         rename => { "response" => "response_new" }
15         gsub => ["referrer", "\\\"", "\""]
16         split => ["clientip", "."]
17     }
18 }
19 output {
20     stdout {
21         codec => "rubydebug"
22     }

```

#20.4 Logstash输出插件（output）

<https://www.elastic.co/guide/en/logstash/current/output-plugins.html>

output是Logstash的最后阶段，一个事件可以经过多个输出，而一旦所有输出处理完成，整个事件就执行完成。一些常用的输出包括：

- file：表示将日志数据写入磁盘上的文件。
- elasticsearch：表示将日志数据发送给Elasticsearch。Elasticsearch可以高效方便和易于查询的保存数据。

1、输出到标准输出(stdout)

```

1  output {
2      stdout {
3          codec => rubydebug
4      }
5  }

```

2、保存为文件（file）

```

1  output {
2      file {
3          path => "/data/log/%{+yyyy-MM-dd}/%{host}_%{+HH}.log"
4      }
5  }

```

3、输出到elasticsearch

```

1  output {
2      elasticsearch {
3          host => ["192.168.1.1:9200", "172.16.213.77:9200"]
4          index => "logstash-%{+YYYY.MM.dd}"
5      }
6  }

```

- host：是一个数组类型的值，后面跟的值是elasticsearch节点的地址与端口，默认端口是9200。可添加多个地址。
- index：写入elasticsearch的索引的名称，这里可以使用变量。Logstash提供了%{+YYYY.MM.dd}这种写法。在语法解析的时候，看到以+号开头的，就会自动认为后面是时间格式，尝试用时间格式来解析后续字符串。这种以天为单位分割的写法，可以很容易的删除老的数据或者搜索指定时间范围内的数据。此外，注意索引名中不能有大写字母。
- manage_template:用来设置是否开启logstash自动管理模板功能，如果设置为false将关闭自动管理模板功能。如果我们自定义了模板，那么应该设置为false。
- template_name:这个配置项用来设置在Elasticsearch中模板的名称。

#20.5 综合案例

```
1  input {
2      file {
3          path => ["D:/ES/logstash-7.3.0/nginx.log"]
4          start_position => "beginning"
5      }
6  }
7
8  filter {
9      grok {
10         match => { "message" => "%{IP:clientip}\ \[%{HTTPDATE:timestamp}\]\ %
11         {QS:referrer}\ %{NUMBER:response}\ %{NUMBER:bytes}" }
12         remove_field => [ "message" ]
13     }
14     date {
15         match => ["timestamp", "dd/MMM/yyyy:HH:mm:ss Z"]
16     }
17     mutate {
18         rename => { "response" => "response_new" }
19         convert => [ "response", "float" ]
20         gsub => ["referrer", "\\", "\\\\" ]
21         remove_field => ["timestamp"]
22         split => ["clientip", "."]
23     }
24 }
25 output {
26     stdout {
27         codec => "rubydebug"
28     }
29 }
30 elasticsearch {
31     host => ["localhost:9200"]
32     index => "logstash-%{+YYYY.MM.dd}"
33 }
34
35 }
```

#21. kibana学习

#21.1基本查询

1是什么：elk中数据展现工具。

2下载：<https://www.elastic.co/cn/downloads/kibana>

3使用：建立索引模式，index partten

discover 中使用DSL搜索。

#21.2 可视化

绘制图形

#21.3 仪表盘

将各种可视化图形放入，形成大屏幕。

#21.4 使用模板数据指导绘图

点击主页的添加模板数据，可以看到很多模板数据以及绘图。

#21.5 其他功能

监控，日志，APM等功能非常丰富。

#22. 集群部署

见部署图

#结点的三个角色

主结点：master节点主要用于集群的管理及索引 比如新增结点、分片分配、索引的新增和删除等。数据结点：data 节点上保存了数据分片，它负责索引和搜索操作。客户端结点：client 节点仅作为请求客户端存在，client的作用也作为负载均衡器，client 节点不存数据，只是将请求均衡转发到其它结点。

通过下边两项参数来配置结点的功能：

```
1 node.master: #是否允许为主结点
2 node.data: #允许存储数据作为数据结点
3 node.ingest: #是否允许成为协调节点
```

四种组合方式：

```
1 master=true, data=true: 即是主结点又是数据结点
2 master=false, data=true: 仅是数据结点
3 master=true, data=false: 仅是主结点，不存储数据
4 master=false, data=false: 即不是主结点也不是数据结点，此时可设置ingest为true表示它是一个客户端。
```

#23. 项目实战

#23.1项目一：ELK用于日志分析

需求：集中收集分布式服务的日志

1逻辑模块程序随时输出日志

```
1 package com.ydlclass.es;
2
```

```

3  import org.junit.Test;
4  import org.junit.runner.RunWith;
5  import org.slf4j.Logger;
6  import org.slf4j.LoggerFactory;
7  import org.springframework.boot.test.context.SpringBootTest;
8  import org.springframework.test.context.junit4.SpringRunner;
9
10 import java.util.Random;
11
12 /**
13  * creste by ydlclass.ydl
14  */
15 @SpringBootTest
16 @RunWith(SpringRunner.class)
17 public class TestLog {
18     private static final Logger LOGGER= LoggerFactory.getLogger(TestLog.class);
19
20     @Test
21     public void testLog(){
22         Random random =new Random();
23
24         while (true){
25             int userid=random.nextInt(10);
26             LOGGER.info("userId:{},send:{})",userid,"hello world.I am "+userid);
27             try {
28                 Thread.sleep(500);
29             } catch (InterruptedException e) {
30                 e.printStackTrace();
31             }
32         }
33     }
34 }

```

2logstash收集日志到es

grok 内置类型

```

1  USERNAME [a-zA-Z0-9._-]+
2  USER %{USERNAME}
3  INT (?:[+-]?(?:[0-9]+))
4  BASE10NUM (?<![0-9.-+])(?>[+-]?(?:[0-9]+(?:\.[0-9]+)?|(?:\.[0-9]+)))
5  NUMBER (?:%{BASE10NUM})
6  BASE16NUM (?<![0-9A-Fa-f])(?:[+-]?(?:0x)?(?:[0-9A-Fa-f]+))
7  BASE16FLOAT \b(?<![0-9A-Fa-f.])?(?:[+-]?(?:0x)?(?:[0-9A-Fa-f]+(?:\.[0-9A-Fa-f]*)?|(?:\.[0-9A-Fa-f]+)))\b
8
9  POSINT \b(?:[1-9][0-9]*)\b
10 NONNEGINT \b(?:[0-9]+)\b
11 WORD \b\w+\b
12 NOTSPACE \S+
13 SPACE \s*
14 DATA .*?
15 GREEDYDATA .*
16 QUOTEDSTRING (?>(?!\\)(?>"(?:\\.|[^\\""]+)"|'(?>\\.|[^\\"']+)'|`(?>\\.|[^\\"`']+)`|(?>`|'|"|`))
17 UUID [A-Fa-f0-9]{8}-(?:[A-Fa-f0-9]{4})-(?:[A-Fa-f0-9]{4})-(?:[A-Fa-f0-9]{12})
18
19 # Networking

```

```

20  MAC (?:%{CISCOMAC}|%{WINDOWSMAC}|%{COMMONMAC})
21  CISCOMAC (?:([A-Fa-f0-9]{4}\.){2}[A-Fa-f0-9]{4})
22  WINDOWSMAC (?:([A-Fa-f0-9]{2}-){5}[A-Fa-f0-9]{2})
23  COMMONMAC (?:([A-Fa-f0-9]{2}:){5}[A-Fa-f0-9]{2})
24  IPV6 ((([0-9A-Fa-f]{1,4}:){7}([0-9A-Fa-f]{1,4}|:))|(([0-9A-Fa-f]{1,4}:){6}(:[0-9A-Fa-f]{1,4}|((25[0-5]|2[0-4]\d|1\d\d|1-9)?\d)(\. (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d))){3}|:))|(([0-9A-Fa-f]{1,4}:){5}((:[0-9A-Fa-f]{1,4}){1,2})|:((25[0-5]|2[0-4]\d|1\d\d|1-9)?\d)(\. (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d))){3}|:)|(([0-9A-Fa-f]{1,4}:){4}(((:[0-9A-Fa-f]{1,4}){1,3})|((:[0-9A-Fa-f]{1,4})?:((25[0-5]|2[0-4]\d|1\d\d|1-9)?\d)(\. (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d))){3}))|:)|(([0-9A-Fa-f]{1,4}:){3}(((:[0-9A-Fa-f]{1,4}){1,4})|((:[0-9A-Fa-f]{1,4}){0,2}:((25[0-5]|2[0-4]\d|1\d\d|1-9)?\d)(\. (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d))){3}))|:)|(([0-9A-Fa-f]{1,4}:){2}(((:[0-9A-Fa-f]{1,4}){1,5})|((:[0-9A-Fa-f]{1,4}){0,3}:((25[0-5]|2[0-4]\d|1\d\d|1-9)?\d)(\. (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d))){3}))|:)|(([0-9A-Fa-f]{1,4}:){1}(((:[0-9A-Fa-f]{1,4}){1,6})|((:[0-9A-Fa-f]{1,4}){0,4}:((25[0-5]|2[0-4]\d|1\d\d|1-9)?\d)(\. (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d))){3}))|:)|(:((:[0-9A-Fa-f]{1,4}){1,7})|((:[0-9A-Fa-f]{1,4}){0,5}:((25[0-5]|2[0-4]\d|1\d\d|1-9)?\d)(\. (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d))){3}))|:)))(%.+)?
25  IPV4 (?<![0-9])(?:([0-9]{1,3}|[0-9]{1,2})[.](?:[0-9]{1,3}|[0-9]{1,2})[.](?:[0-9]{1,3}|[0-9]{1,2})[.](?:[0-9]{1,3}|[0-9]{1,2}))?![0-9])
26  IP (?:%{IPV6}|%{IPV4})
27  HOSTNAME \b(?:[0-9A-Za-z][0-9A-Za-z-]{0,62})(?:\.(?:[0-9A-Za-z][0-9A-Za-z-]{0,62}))*(\.|\b)
28  HOST %{HOSTNAME}
29  IPORHOST (?:%{HOSTNAME}|%{IP})
30  HOSTPORT %{IPORHOST}:%{POSINT}
31
32  # paths
33  PATH (?:%{UNIXPATH}|%{WINPATH})
34  UNIXPATH (?>/(?>[\w_!$@:.,-]+|\\\.)*+
35  TTY (?:/dev/(pts|tty([pq])?) (\w+)?/?(?:[0-9]+))
36  WINPATH (?>[A-Za-z]+:|\\)(?:\\[\^\\?]*)*+
37  URIPROTO [A-Za-z]+(\+[A-Za-z]+)?
38  URIHOST %{IPORHOST}(?::%{POSINT:port})?
39  # uripath comes loosely from RFC1738, but mostly from what Firefox
40  # doesn't turn into %XX
41  URIPATH (?:/[A-Za-z0-9$.+!*'(){}~,;:@#%_\-]*)+
42  #URIPARAM \?(?:[A-Za-z0-9]+(?:=(?:[^\&]*)?(?:&(?:[A-Za-z0-9]+(?:=(?:[^\&]*)?)?)*)?)?
43  URIPARAM \?[A-Za-z0-9$.+!*'(){}~,;:@#%&/=-;_?\\-[\]]*
44  URIPATHPARAM %{URIPATH}(?:%{URIPARAM})?
45  URI %{URIPROTO}://(?:%{USER}(?:[^\@]*)?@)(?:%{URIHOST})(?:%{URIPATHPARAM})?
46
47  # Months: January, Feb, 3, 03, 12, December
48  MONTH \b(?:Jan(?:uary)?|Feb(?:ruary)?|Mar(?:ch)?|Apr(?:il)?|May|Jun(?:e)?|Jul(?:y)?|Aug(?:ust)?|Sep(?:tember)?|Oct(?:ober)?|Nov(?:ember)?|Dec(?:ember)?)\b
49  MONTHNUM (?:0?[1-9]|1[0-2])
50  MONTHNUM2 (?:0[1-9]|1[0-2])
51  MONTHDAY (?:([0-9]|([12][0-9])|([3][01]))|([1-9]))
52
53  # Days: Monday, Tue, Thu, etc...
54  DAY (?:Mon(?:day)?|Tue(?:sday)?|Wed(?:nesday)?|Thu(?:rsday)?|Fri(?:day)?|Sat(?:urday)?|Sun(?:day)?)
55
56  # Years?
57  YEAR (?>\d\d){1,2}
58  HOUR (?:2[0123]|01)?[0-9])

```

```

59  MINUTE (?:[0-5][0-9])
60  # '60' is a leap second in most time standards and thus is valid.
61  SECOND (?:(?:[0-5]?[0-9]|60)(?:[:.,][0-9]+)?)
62  TIME (?!<[0-9])%{HOUR}:%{MINUTE}(?::%{SECOND})(?![0-9])
63  # timestamp is YYYY/MM/DD-HH:MM:SS.UUUU (or something like it)
64  DATE_US %{MONTHNUM}[/-]%{MONTHDAY}[/-]%{YEAR}
65  DATE_EU %{MONTHDAY}[./-]%{MONTHNUM}[./-]%{YEAR}
66  ISO8601_TIMEZONE (?:Z|[-+]%{HOUR}(?::%{MINUTE}))
67  ISO8601_SECOND (?::%{SECOND}|60)
68  TIMESTAMP_ISO8601 %{YEAR}-%{MONTHNUM}-%{MONTHDAY}[T ]%{HOUR}:%{MINUTE}(?::%{SECOND})?%{ISO8601_TIMEZONE}?
69  DATE %{DATE_US}|%{DATE_EU}
70  DATESTAMP %{DATE}[- ]%{TIME}
71  TZ (?:[PMCE][SD]T|UTC)
72  DATESTAMP_RFC822 %{DAY} %{MONTH} %{MONTHDAY} %{YEAR} %{TIME} %{TZ}
73  DATESTAMP_RFC2822 %{DAY}, %{MONTHDAY} %{MONTH} %{YEAR} %{TIME} %
    {ISO8601_TIMEZONE}
74  DATESTAMP_OTHER %{DAY} %{MONTH} %{MONTHDAY} %{TIME} %{TZ} %{YEAR}
75  DATESTAMP_EVENTLOG %{YEAR}%{MONTHNUM2}%{MONTHDAY}%{HOUR}%{MINUTE}%{SECOND}
76
77  # Syslog Dates: Month Day HH:MM:SS
78  SYSLOGTIMESTAMP %{MONTH} +%{MONTHDAY} %{TIME}
79  PROG (?:[\\w._/-]+)
80  SYSLOGPROG %{PROG:program}(?:\\[%{POSINT:pid}\\])?
81  SYSLOGHOST %{IPORHOST}
82  SYSLOGFACILITY <%{NONNEGINT:facility}%.%{NONNEGINT:priority}>
83  HTTPDATE %{MONTHDAY}/%{MONTH}/%{YEAR}:%{TIME} %{INT}
84
85  # Shortcuts
86  QS %{QUOTEDSTRING}
87
88  # Log formats
89  SYSLOGBASE %{SYSLOGTIMESTAMP:timestamp} (?::%{SYSLOGFACILITY} )?%
    {SYSLOGHOST:logsource} %{SYSLOGPROG}:
90  COMMONAPACHELOG %{IPORHOST:clientip} %{USER:ident} %{USER:auth} \\[%
    {HTTPDATE:timestamp}\\] "(?::%{WORD:verb} %{NOTSPACE:request}(?: HTTP/%
    {NUMBER:httpversion})?|%{DATA:rawrequest})" %{NUMBER:response} (?:%
    {NUMBER:bytes}|-)
91  COMBINEDAPACHELOG %{COMMONAPACHELOG} %{QS:referrer} %{QS:agent}
92
93  # Log Levels
94  LOGLEVEL ([Aa]lert|ALERT|[Tt]race|TRACE|[Dd]ebug|DEBUG|[Nn]otice|NOTICE|
    [Ii]nfo|INFO|[Ww]arn?(?:ing)?|WARN?(?:ING)?|[Ee]rr?(?:or)?|ERR?(?:OR)?|[Cc]rit?
    (?::ical)?|CRIT?(?::ICAL)?|[Ff]atal|FATAL|[Ss]evere|SEVERE|EMERG(?::ENCY)?|
    [Ee]merg(?::ency)?)

```

写logstash配置文件。

3kibana展现数据

#23.2项目二：学成在线站内搜索模块

#1mysql导入course_pub表

#2创建索引xc_course

#3创建映射

```
1  PUT /xc_course
2  {
3    "settings": {
4      "number_of_shards": 1,
5      "number_of_replicas": 0
6    },
7    "mappings": {
8      "properties": {
9        "description" : {
10          "analyzer" : "ik_max_word",
11          "search_analyzer": "ik_smart",
12          "type" : "text"
13        },
14        "grade" : {
15          "type" : "keyword"
16        },
17        "id" : {
18          "type" : "keyword"
19        },
20        "mt" : {
21          "type" : "keyword"
22        },
23        "name" : {
24          "analyzer" : "ik_max_word",
25          "search_analyzer": "ik_smart",
26          "type" : "text"
27        },
28        "users" : {
29          "index" : false,
30          "type" : "text"
31        },
32        "charge" : {
33          "type" : "keyword"
34        },
35        "valid" : {
36          "type" : "keyword"
37        },
38        "pic" : {
39          "index" : false,
40          "type" : "keyword"
41        },
42        "qq" : {
43          "index" : false,
44          "type" : "keyword"
45        },
46        "price" : {
47          "type" : "float"
48        },
49        "price_old" : {
50          "type" : "float"
```

```

51         },
52         "st" : {
53             "type" : "keyword"
54         },
55         "status" : {
56             "type" : "keyword"
57         },
58         "studymodel" : {
59             "type" : "keyword"
60         },
61         "teachmode" : {
62             "type" : "keyword"
63         },
64         "teachplan" : {
65             "analyzer" : "ik_max_word",
66             "search_analyzer": "ik_smart",
67             "type" : "text"
68         },
69         "expires" : {
70             "type" : "date",
71             "format": "yyyy-MM-dd HH:mm:ss"
72         },
73         "pub_time" : {
74             "type" : "date",
75             "format": "yyyy-MM-dd HH:mm:ss"
76         },
77         "start_time" : {
78             "type" : "date",
79             "format": "yyyy-MM-dd HH:mm:ss"
80         },
81         "end_time" : {
82             "type" : "date",
83             "format": "yyyy-MM-dd HH:mm:ss"
84         }
85     }
86 }
87 }

```

#4logstash创建模板文件

Logstash的工作是从MySQL中读取数据，向ES中创建索引，这里需要提前创建mapping的模板文件以便logstash使用。

在logstash的config目录创建xc_course_template.json，内容如下：

```

1  {
2      "mappings" : {
3          "doc" : {
4              "properties" : {
5                  "charge" : {
6                      "type" : "keyword"
7                  },
8                  "description" : {
9                      "analyzer" : "ik_max_word",
10                     "search_analyzer" : "ik_smart",
11                     "type" : "text"
12                 },

```

```
13     "end_time" : {
14         "format" : "yyyy-MM-dd HH:mm:ss",
15         "type" : "date"
16     },
17     "expires" : {
18         "format" : "yyyy-MM-dd HH:mm:ss",
19         "type" : "date"
20     },
21     "grade" : {
22         "type" : "keyword"
23     },
24     "id" : {
25         "type" : "keyword"
26     },
27     "mt" : {
28         "type" : "keyword"
29     },
30     "name" : {
31         "analyzer" : "ik_max_word",
32         "search_analyzer" : "ik_smart",
33         "type" : "text"
34     },
35     "pic" : {
36         "index" : false,
37         "type" : "keyword"
38     },
39     "price" : {
40         "type" : "float"
41     },
42     "price_old" : {
43         "type" : "float"
44     },
45     "pub_time" : {
46         "format" : "yyyy-MM-dd HH:mm:ss",
47         "type" : "date"
48     },
49     "qq" : {
50         "index" : false,
51         "type" : "keyword"
52     },
53     "st" : {
54         "type" : "keyword"
55     },
56     "start_time" : {
57         "format" : "yyyy-MM-dd HH:mm:ss",
58         "type" : "date"
59     },
60     "status" : {
61         "type" : "keyword"
62     },
63     "studymodel" : {
64         "type" : "keyword"
65     },
66     "teachmode" : {
67         "type" : "keyword"
68     },
69     "teachplan" : {
70         "analyzer" : "ik_max_word",
```



```

71         "search_analyzer" : "ik_smart",
72         "type" : "text"
73     },
74     "users" : {
75         "index" : false,
76         "type" : "text"
77     },
78     "valid" : {
79         "type" : "keyword"
80     }
81 }
82 }
83 },
84 "template" : "xc_course"
85 }

```

#5logstash配置mysql.conf

1、ES采用UTC时区问题

ES采用UTC 时区，比北京时间早8小时，所以ES读取数据时让最后更新时间加8小时

where timestamp > date_add(:sql_last_value,INTERVAL 8 HOUR)

2、logstash每个执行完成会在/config/logstash_metadata记录执行时间下次以此时间为基准进行增量同步数据到索引库。

#6启动

```
1 .\logstash.bat -f ..\config\mysql.conf
```

7后端代码

7.1Controller

```

1  @RestController
2  @RequestMapping("/search/course")
3  public class EsCourseController {
4      @Autowired
5      EsCourseService esCourseService;
6
7      @GetMapping(value="/list/{page}/{size}")
8      public QueryResponseResult<CoursePub> list(@PathVariable("page") int page,
9          @PathVariable("size") int size, CourseSearchParams courseSearchParams) {
10         return esCourseService.list(page, size, courseSearchParams);
11     }
12 }

```

7.2

```

1  @Service
2  public class EsCourseService {
3      @Value("${heima.course.source_field}")
4      private String source_field;
5
6      @Autowired
7      RestHighLevelClient restHighLevelClient;

```

```

8
9      //课程搜索
10     public QueryResponseResult<CoursePub> list(int page, int size,
CourseSearchParam courseSearchParam) {
11         if (courseSearchParam == null) {
12             courseSearchParam = new CourseSearchParam();
13         }
14         //1创建搜索请求对象
15         SearchRequest searchRequest = new SearchRequest("xc_course");
16
17         SearchSourceBuilder searchSourceBuilder = new SearchSourceBuilder();
18         //过滤源字段
19         String[] source_field_array = source_field.split(",");
20         searchSourceBuilder.fetchSource(source_field_array, new String[]{});
21         //创建布尔查询对象
22         BoolQueryBuilder boolQueryBuilder = QueryBuilders.boolQuery();
23         //搜索条件
24         //根据关键字搜索
25         if (StringUtils.isEmpty(courseSearchParam.getKeyword())) {
26             MultiMatchQueryBuilder multiMatchQueryBuilder =
QueryBuilders.multiMatchQuery(courseSearchParam.getKeyword(), "name",
"description", "teachplan")
27                 .minimumShouldMatch("70%")
28                 .field("name", 10);
29             boolQueryBuilder.must(multiMatchQueryBuilder);
30         }
31         if (StringUtils.isEmpty(courseSearchParam.getMt())) {
32             //根据一级分类
33             boolQueryBuilder.filter(QueryBuilders.termQuery("mt",
courseSearchParam.getMt()));
34         }
35         if (StringUtils.isEmpty(courseSearchParam.getSt())) {
36             //根据二级分类
37             boolQueryBuilder.filter(QueryBuilders.termQuery("st",
courseSearchParam.getSt()));
38         }
39         if (StringUtils.isEmpty(courseSearchParam.getGrade())) {
40             //根据难度等级
41             boolQueryBuilder.filter(QueryBuilders.termQuery("grade",
courseSearchParam.getGrade()));
42         }
43
44         //设置boolQueryBuilder到searchSourceBuilder
45         searchSourceBuilder.query(boolQueryBuilder);
46         //设置分页参数
47         if (page <= 0) {
48             page = 1;
49         }
50         if (size <= 0) {
51             size = 12;
52         }
53         //起始记录下标
54         int from = (page - 1) * size;
55         searchSourceBuilder.from(from);
56         searchSourceBuilder.size(size);
57
58         //设置高亮
59         HighlightBuilder highlightBuilder = new HighlightBuilder();

```

```

60     highlightBuilder.preTags("<font class='eslight'>");
61     highlightBuilder.postTags("</font>");
62     //设置高亮字段
63     //     <font class='eslight'>node</font>学习
64     highlightBuilder.fields().add(new HighlightBuilder.Field("name"));
65     searchSourceBuilder.highlighter(highlightBuilder);
66
67     searchRequest.source(searchSourceBuilder);
68
69     QueryResult<CoursePub> queryResult = new QueryResult();
70     List<CoursePub> list = new ArrayList<CoursePub>();
71     try {
72         //2执行搜索
73         SearchResponse searchResponse =
74         restHighLevelClient.search(searchRequest, RequestOptions.DEFAULT);
75         //3获取响应结果
76         SearchHits hits = searchResponse.getHits();
77         long totalHits=hits.getTotalHits().value;
78         //匹配的总记录数
79         //     long totalHits = hits.totalHits;
80         queryResult.setTotal(totalHits);
81         SearchHit[] searchHits = hits.getHits();
82         for (SearchHit hit : searchHits) {
83             CoursePub coursePub = new CoursePub();
84             //源文档
85             Map<String, Object> sourceAsMap = hit.getSourceAsMap();
86             //取出id
87             String id = (String) sourceAsMap.get("id");
88             coursePub.setId(id);
89             //取出name
90             String name = (String) sourceAsMap.get("name");
91             //取出高亮字段name
92             Map<String, HighlightField> highlightFields =
93             hit.getHighlightFields();
94             if (highlightFields != null) {
95                 HighlightField highlightFieldName =
96                 highlightFields.get("name");
97                 if (highlightFieldName != null) {
98                     Text[] fragments = highlightFieldName.fragments();
99                     StringBuffer stringBuffer = new StringBuffer();
100                     for (Text text : fragments) {
101                         stringBuffer.append(text);
102                     }
103                     name = stringBuffer.toString();
104                 }
105             }
106             coursePub.setName(name);
107             //图片
108             String pic = (String) sourceAsMap.get("pic");
109             coursePub.setPic(pic);
110             //价格
111             Double price = null;
112             try {
113                 if (sourceAsMap.get("price") != null) {
114                     price = (Double) sourceAsMap.get("price");
115                 }
116             } catch (Exception e) {

```

```
115         e.printStackTrace();
116     }
117     coursePub.setPrice(price);
118     //旧价格
119     Double price_old = null;
120     try {
121         if (sourceAsMap.get("price_old") != null) {
122             price_old = (Double) sourceAsMap.get("price_old");
123         }
124     } catch (Exception e) {
125         e.printStackTrace();
126     }
127     coursePub.setPrice_old(price_old);
128     //将coursePub对象放入list
129     list.add(coursePub);
130 }
131
132
133     } catch (IOException e) {
134         e.printStackTrace();
135     }
136
137     queryResult.setList(list);
138     QueryResponseResult<CoursePub> queryResponseResult = new
QueryResponseResult<CoursePub>(CommonCode.SUCCESS, queryResult);
139
140     return queryResponseResult;
141 }
142
143
144 }
```