

管理信息系统 NoSQL 实验报告

180812-17374463-王凯东

实验目的：

在本地计算机上使用 Python 完成示例数据集（news.txt）上 ES 和 MySQL 数据写入、查询过程，对比两种数据库的功能和性能差异，程序输出和效率。

1、安装 Java1.8、Python3.7.0 并成功配置环境变量

```
C:\WINDOWS\system32\cmd.exe
Microsoft Windows [版本 10.0.18363.719]
(c) 2019 Microsoft Corporation. 保留所有权利。

C:\Users\Windows10>java -version
java version "1.8.0_241"
Java(TM) SE Runtime Environment (build 1.8.0_241-b07)
Java HotSpot(TM) 64-Bit Server VM (build 25.241-b07, mixed mode)

C:\Users\Windows10>pip

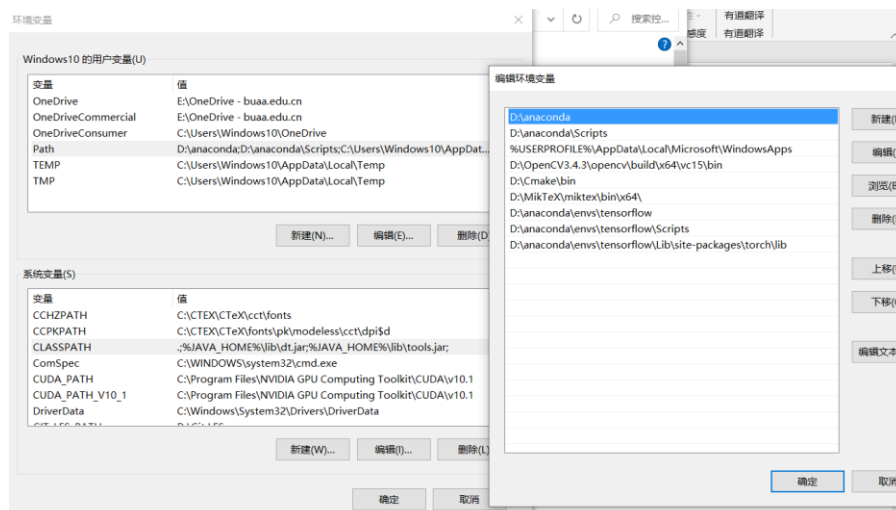
Usage:
  pip <command> [options]

Commands:
  install             Install packages.
  download            Download packages.
  uninstall           Uninstall packages.
  freeze              Output installed packages in requirements format.
  list                List installed packages.
  show                Show information about installed packages.
  check               Verify installed packages have compatible dependencies.
  config              Manage local and global configuration.
  search              Search PyPI for packages.
  wheel               Build wheels from your requirements.
  hash                Compute hashes of package archives.
  completion          A helper command used for command completion.
  debug               Show information useful for debugging.
  help                Show help for commands.

General Options:
```

```
C:\WINDOWS\system32\cmd.exe - python
Microsoft Windows [版本 10.0.18363.719]
(c) 2019 Microsoft Corporation. 保留所有权利。

C:\Users\Windows10>python
Python 3.7.0 (default, Jun 28 2018, 08:04:48) [MSC v.1912 64 bit (AMD64)] :: Anaconda, Inc. on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>> from elasticsearch import elasticsearch
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
ImportError: cannot import name 'elasticsearch' from 'elasticsearch' (D:\anaconda\lib\site-packages\elasticsearch\_init_.py)
>>> from elasticsearch import Elasticsearch
>>> import pymysql
>>>
```

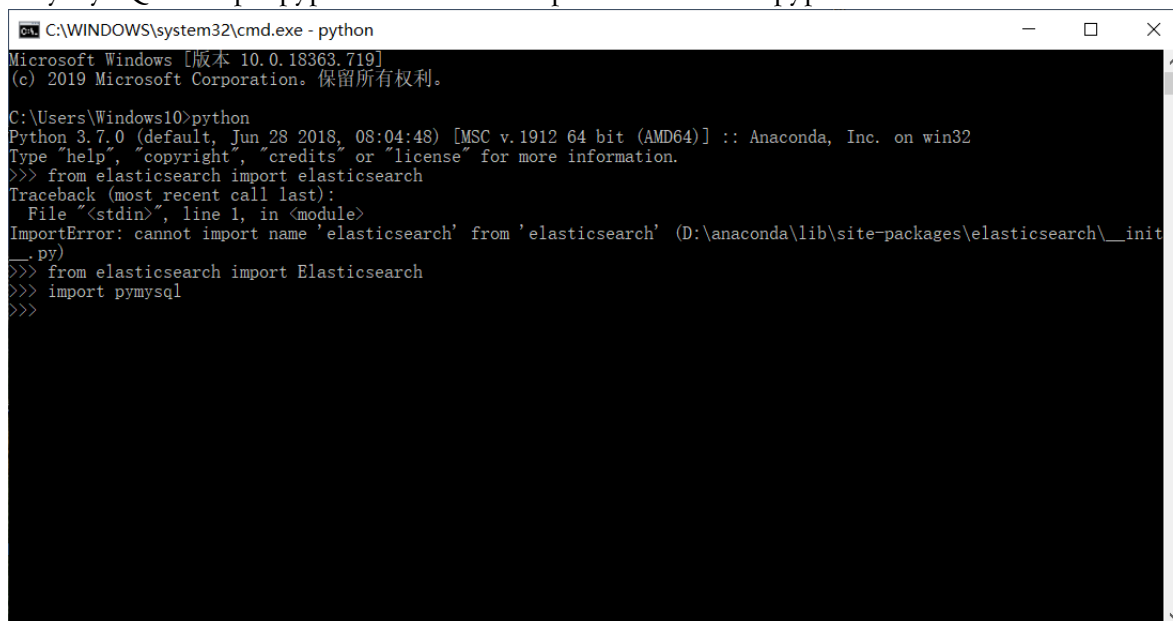


2、安装 elasticsearch 库和 pymysql 库

在命令行执行：

```
pip install elasticsearch==5.4.0 -i http://pypi.douban.com/simple --trusted-host pypi.douban.com;
```

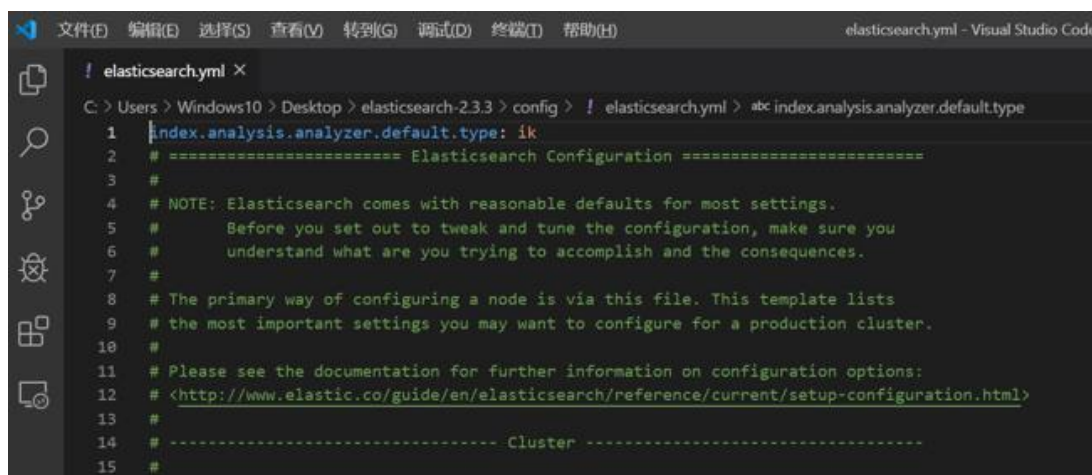
```
pip install PyMySQL -i http://pypi.douban.com/simple --trusted-host pypi.douban.com
```



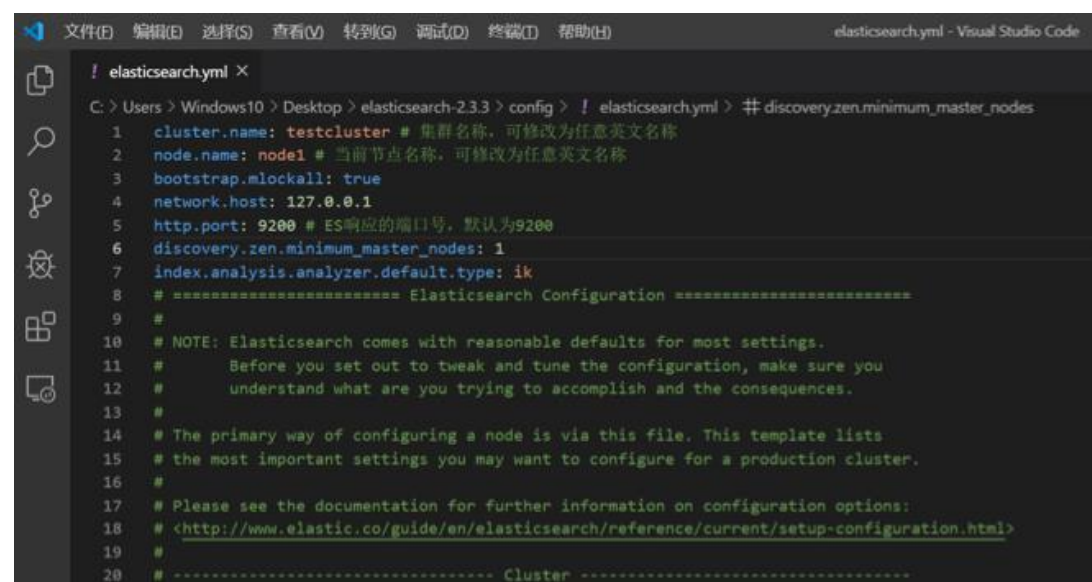
```
C:\WINDOWS\system32\cmd.exe - python
Microsoft Windows [版本 10.0.18363.719]
(c) 2019 Microsoft Corporation. 保留所有权利。

C:\Users\Windows10>python
Python 3.7.0 (default, Jun 28 2018, 08:04:48) [MSC v.1912 64 bit (AMD64)] :: Anaconda, Inc. on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>> from elasticsearch import elasticsearch
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
ImportError: cannot import name 'elasticsearch' from 'elasticsearch' (D:\anaconda\lib\site-packages\elasticsearch\__init__
.py)
>>> from elasticsearch import Elasticsearch
>>> import pymysql
>>>
```

3、安装 IK 分词插件并创建结点 node1



```
elasticsearch.yml - Visual Studio Code
! elasticsearch.yml X
C:\Users\Windows10\Desktop> elasticsearch-2.3.3> config> ! elasticsearch.yml> # index.analysis.analyzer.default.type
1 index.analysis.analyzer.default.type: ik
2 # ===== Elasticsearch Configuration =====
3 #
4 # NOTE: Elasticsearch comes with reasonable defaults for most settings.
5 # Before you set out to tweak and tune the configuration, make sure you
6 # understand what are you trying to accomplish and the consequences.
7 #
8 # The primary way of configuring a node is via this file. This template lists
9 # the most important settings you may want to configure for a production cluster.
10 #
11 # Please see the documentation for further information on configuration options:
12 # <http://www.elastic.co/guide/en/elasticsearch/reference/current/setup-configuration.html>
13 #
14 # ----- Cluster -----
15 #
```



```
elasticsearch.yml - Visual Studio Code
! elasticsearch.yml X
C:\Users\Windows10\Desktop> elasticsearch-2.3.3> config> ! elasticsearch.yml> # discovery.zen.minimum_master_nodes
1 cluster.name: testcluster # 集群名称，可修改为任意英文名称
2 node.name: node1 # 当前节点名称，可修改为任意英文名称
3 bootstrap.mlockall: true
4 network.host: 127.0.0.1
5 http.port: 9200 # ES响应的端口号，默认为9200
6 discovery.zen.minimum_master_nodes: 1
7 index.analysis.analyzer.default.type: ik
8 # ===== Elasticsearch Configuration =====
9 #
10 # NOTE: Elasticsearch comes with reasonable defaults for most settings.
11 # Before you set out to tweak and tune the configuration, make sure you
12 # understand what are you trying to accomplish and the consequences.
13 #
14 # The primary way of configuring a node is via this file. This template lists
15 # the most important settings you may want to configure for a production cluster.
16 #
17 # Please see the documentation for further information on configuration options:
18 # <http://www.elastic.co/guide/en/elasticsearch/reference/current/setup-configuration.html>
19 #
20 # ----- Cluster -----
```

5、构建第二个结点 node2，再次刷新



7、分别创建 ES 和 MySQL 数据库和表结构

```
MySQL 5.7 Command Line Client
Enter password: *****
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 4
Server version: 5.7.29-log MySQL Community Server (GPL)

Copyright (c) 2000, 2020, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> create database test;
Query OK, 1 row affected (0.00 sec)

mysql> use test;
Database changed
mysql> CREATE TABLE news (
->   id varchar(255) NOT NULL,
->   url varchar(255) NOT NULL,
->   datestr varchar(255) NOT NULL,
->   title varchar(255) NOT NULL,
->   content varchar(255) NOT NULL,
->   PRIMARY KEY (id)
-> ) ENGINE=InnoDB DEFAULT CHARSET=utf8;
Query OK, 0 rows affected (0.05 sec)

mysql>
```

```
C: > Users > Windows10 > Desktop > code > es_create.py > ...

4   from config import ES_INDEX_NAME, ES_INDEX_DOC_TYPE, es
5
6   def create_mapping():
7       index_info = {
8           'settings':{
9               'number_of_replicas': 1,
10              'number_of_shards': 5
11          },
12          'mappings':{
13              ES_INDEX_DOC_TYPE:{
14                  'properties':{
15                      "url":{
16                          'type': 'string',
17                          'index': 'not_analyzed'
18                      },
19                      "title":{
20                          'type': 'string',
21                          'index': 'analyzed'
22                      },
23                      "content":{
24                          'type': 'string',
25                          'index': 'analyzed'
26                      },
27                      "date":{
28                          'type': 'string',
29                          'index': 'not_analyzed'
30                      }
31                  }
32              }
33          }
34      }
35
```

• 使用 ES 写入:

信息 ▼ 动作 ▼

[illegible]

• 使用 MySQL 写入:

```
10 # read data from txt
11 data_list = []
12 with open("news.txt", "r", encoding='UTF-8') as f:
13     for line in f:
14         url, date, title, content = line.strip().split("|.")
15         data_list.append((url, url, title, content, date))
16
17 # Connect to the database
18 connection = pymysql.connect(host=MYSQL_HOST,
19                               user=MYSQL_USERNAME,
20                               password=MYSQL_PASSWORD,
21                               db=MYSQL_DB,
22                               charset='utf8mb4',
23                               cursorclass=pymysql.cursors.DictCursor)
24
25 # write data
26 count = 0
27 tb = time.time()
28 with connection.cursor() as cursor:
29     for data in data_list:
30         # Create a new record
31         sql = "INSERT INTO `" + MYSQL_TABLE + "` (`id`, `url`, `title`, `content`, `datestr`) VALUES (%s, %s, %s, %s, %s)"
32         cursor.execute(sql, data)
33
34         # connection is not autocommit by default. So you must commit to save your changes.
35         connection.commit()
36         count += 1
37
38 print ("MySQL writing %s items need %s seconds" % (count, time.time() - tb))
```

• 使用 ES 和 MySQL 进行写入分别耗时 17.2928s 和 1.8589s。

```
C:\WINDOWS\system32\cmd.exe
Microsoft Windows [版本 10.0.18363.719]
(c) 2019 Microsoft Corporation。保留所有权利。

C:\Users\Windows10>cd C:\Users\Windows10\Desktop\code

C:\Users\Windows10\Desktop\code>python es_write.py
ES writing 630 items need 17.292790174484253 seconds

C:\Users\Windows10\Desktop\code>python mysql_write.py
Traceback (most recent call last):
  File "mysql_write.py", line 32, in <module>
    cursor.execute(sql, data)
  File "D:\anaconda\lib\site-packages\pymysql\cursors.py", line 170, in execute
    result = self._query(query)
  File "D:\anaconda\lib\site-packages\pymysql\cursors.py", line 328, in _query
    conn.query(q)
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 517, in query
    self._affected_rows = self._read_query_result(unbuffered=unbuffered)
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 732, in _read_query_result
    result.read()
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 1075, in read
    first_packet = self.connection._read_packet()
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 684, in _read_packet
    packet.check_error()
  File "D:\anaconda\lib\site-packages\pymysql\protocol.py", line 220, in check_error
    err.raise_mysql_exception(self._data)
  File "D:\anaconda\lib\site-packages\pymysql\err.py", line 109, in raise_mysql_exception
    raise errorclass(errno, errval)
pymysql.err.DataError: (1406, "Data too long for column 'content' at row 1")

C:\Users\Windows10\Desktop\code>python mysql_write.py
D:\anaconda\lib\site-packages\pymysql\cursors.py:329: Warning: (1265, "Data truncated for column 'content' at row 1")
  self._do_get_result()
MySQL writing 630 items need 1.8588917255401611 seconds

C:\Users\Windows10\Desktop\code>
```

9、执行全文搜索

- 使用 ES 执行搜索代码：

```
4 import json
5 import time
6 from config import es, ES_INDEX_NAME, ES_INDEX_DOC_TYPE
7
8
9 def query_by_keyword():
10     query_body = {
11         "query": {
12             "match": { "title": "中国" } # 新闻标题中包含中国关键词
13         },
14         "size": 10000
15     }
16
17     tb = time.time()
18     results = es.search(index=ES_INDEX_NAME, doc_type=ES_INDEX_DOC_TYPE, body=query_body)["hits"]["hits"]
19
20     ff = open("search_results.txt", "w", encoding="UTF-8")
21     for item in results:
22         data = item["_source"]
23         ff.write("%s|. %s|. %s|. %s\n" % (data["url"], data["date"], data["title"], data["content"]))
24     ff.close()
25     tb1 = time.time()
26     print ("ES search data time (seconds): ", tb1 - tb)
27
28
29 if __name__ == "__main__":
30     query_by_keyword()
31
```

- 使用 MySQL 执行搜索代码：

```
4 import time
5 import pymysql.cursors
6
7 from config import MYSQL_HOST, MYSQL_USERNAME, MYSQL_PASSWORD, MYSQL_DB, MYSQL_TABLE
8
9
10 # Connect to the database
11 connection = pymysql.connect(host=MYSQL_HOST,
12                               user=MYSQL_USERNAME,
13                               password=MYSQL_PASSWORD,
14                               db=MYSQL_DB,
15                               charset='utf8mb4',
16                               cursorclass=pymysql.cursors.DictCursor)
17
18 tb = time.time()
19
20 with connection.cursor() as cursor:
21     # Read a single record
22     sql = "SELECT * FROM `" + MYSQL_TABLE + "` WHERE `title` like %s"
23     cursor.execute(sql, ("%中国%",))
24     result = cursor.fetchall()
25
26     ff = open("search_results.txt", "w", encoding="UTF-8")
27     for r in result:
28         ff.write("%s|. %s|. %s|. %s\n" % (r["url"], r["datestr"], r["title"], r["content"]))
29     ff.close()
30
31 tb1 = time.time()
32 print ("MySQL search data time (seconds): ", tb1 - tb)
```



search_results - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

http://news.163.com/17/0530/15/CLMOCL3K00018750.html|2017-05-30|[英媒质疑“中国造”蒙内铁路 肯尼亚人:是他们无知
网易新闻]|由中国路桥承建的蒙内铁路，全长约480公里，总投资38亿美元，是肯尼亚独立以来修建的最大基建项目。百年前，英国殖民者花费550万英镑，在东非主持修筑一条铁路，因为修筑艰难等原因，被称为“疯狂铁路”。百年后，中国承建的蒙内铁路与之蒙巴萨-内罗毕段轨铁路基本并行，引发一些英国媒体质疑：这不过是另一条“疯狂铁路”。事实果真如此吗？肯尼亚：铁轨上的国家 “中国造”蒙巴萨-内罗毕标轨铁路即将开通的消息，让肯尼亚这个东非国家兴奋不已。在肯尼亚内罗毕的铁路博物馆里，退役火车头和历史老照片讲述着肯尼亚第一条铁路的故事。“肯尼亚是诞生在铁轨上的国家。”博物馆的小册子如此介绍。米轨铁路的轨距为1米，属于窄轨铁路。百年前，“肯尼亚”这个名字尚不存在，对于英国殖民者而言，这里只是一片无名荒原，横亘在他们和乌干达丰富的矿产资源之间。为打开乌干达大门，加强对“英属东非”的控制，也为了遏制当时已臭名昭著的奴隶贸易，英国政府决定在东非修筑一条铁路，连接印度洋上的蒙巴萨港和内陆的乌干达，并命名其为“乌干达铁路”。对于不了解非洲的英国政客来说，斥巨资在“黑色大陆”上修铁路的决定荒唐可笑。一位英国议员甚至赋诗讽刺其是“疯狂铁路”（lunatic line）。疯狂铁路如今，这条931公里铁路仍拥有此名，原因之一是当年的修建艰难异常，超乎想象。“当时修建铁路的印度劳工没有现在的大型机械，很多铺设工作只能依靠简单的工具，甚至是徒手。”铁路博物馆馆长巴拉萨介绍说。当地部落把铁路称为“铁蛇”，不满“铁蛇”入侵拿起武器反抗，也使得修建举步维艰。根据博物馆的资料显示：从1896年修建开始到1901年铁路竣工，共有2493名铁路工人死亡。他们多是死于当地部落袭击，疟疾、痢疾等传染病，以及令人闻风丧胆的“豪沃的食人狮”。狮子多次袭击铁路工人，制造恐慌，甚至在1900年6月将正在巡查铁路的英国官员拖出车厢咬死。“每一英里米轨铁路，都是由四条铁路工人的生命铺就。”独立之路对后来乘着火车欣赏殖民地草原美景的英国贵妇和探险家们来说，这恐怕只是下午茶时的八卦谈资，当时一幅“乌干达铁路游”广告甚至画着一位火车乘客拍手逗弄草原上的狮子。英国殖民者恐怕不曾料到，这条耗费巨大人力物力修建的铁路，却最终帮助肯尼亚走向独立。“当时的许多爱国人士就是乘着火车，到肯尼亚四处演说，为独立运动争取支持。一些武器也通过铁路运送到了反抗军的手里。”巴拉萨说。一条英国殖民者修建的铁路，却最终帮助终结了殖民统治，这或许是肯尼亚的米轨铁路给历史留下的最大吊诡。同时，铁路所经之处大农场纷纷兴起，肯尼亚的咖啡和牛肉借助铁路运输走向海外，经济的发展使得肯尼亚更有底气对英国统治者说不。铁路的修建还带动了城市的形成。如今的东非大都市内罗毕，便是起源于铁路工人与当地马赛人以物换物的市场。新铁路 新希望而一百年前乘着火车奔向独立的肯尼亚，如今期待另一条铁路助其实现经济腾飞：米轨铁路如今已运营超过百年，自然老化和缺乏维护使得其速度从最初设计的每小时70公里下降到了40公里，肯尼亚大部分的货物运输也因此转移到了公路上。同许多肯尼亚人一样，巴拉萨如今期待设计客运时速120公里（货运80公里）的蒙内铁路能重新激活肯尼亚的铁路系统，并推动肯尼亚的工业化。“未来乘坐蒙内铁路客车可以在一天之内往返内罗毕和蒙巴萨，这是以前坐火车和大巴都无法实现的。“货物能在一天之内从蒙巴萨抵达内罗毕，而靠公路运输则可能要好几天。”巴拉萨说，随着人员与货物流动成本的降低，沿线的商业和工业都将获益。绝对不是“疯狂铁路”而在巴拉萨眼里，米轨铁路曾为肯尼亚带来了了不起的变化，“疯狂铁路”这个称呼只能说明欧洲殖民者对非洲的无知和轻视。修建标轨铁路更是肯尼亚做出的一个“明智选择”，他说，中国公司出色地完成了建造任务，也向世界证明了这一点。“蒙内铁路绝不是另一条疯狂的铁路，这是一条清醒而明智的铁路。”巴拉萨非常肯定地说。

http://news.163.com/17/0223/00/CDTUAFAQ000187V5.html|2017-02-23|13岁雏明将任中国篮协主席? 年轻“掌门”非个例

- 使用 ES 和 MySQL 进行搜索分别耗时 2.2339s 和 0.0314s

```
C:\WINDOWS\system32\cmd.exe
ES writing 630 items need 17.292790174484253 seconds
C:\Users\Windows10\Desktop\code>python mysql_write.py
Traceback (most recent call last):
  File "mysql_write.py", line 32, in <module>
    cursor.execute(sql, data)
  File "D:\anaconda\lib\site-packages\pymysql\cursors.py", line 170, in execute
    result = self._query(query)
  File "D:\anaconda\lib\site-packages\pymysql\cursors.py", line 328, in _query
    conn.query(q)
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 517, in query
    self._affected_rows = self._read_query_result(unbuffered=unbuffered)
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 732, in _read_query_result
    result.read()
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 1075, in read
    first_packet = self.connection._read_packet()
  File "D:\anaconda\lib\site-packages\pymysql\connections.py", line 684, in _read_packet
    packet.check_error()
  File "D:\anaconda\lib\site-packages\pymysql\protocol.py", line 220, in check_error
    err.raise_mysql_exception(self._data)
  File "D:\anaconda\lib\site-packages\pymysql\err.py", line 109, in raise_mysql_exception
    raise errorclass(errno, errval)
pymysql.err.DataError: (1406, "Data too long for column 'content' at row 1")
C:\Users\Windows10\Desktop\code>python mysql_write.py
D:\anaconda\lib\site-packages\pymysql\cursors.py:329: Warning: (1265, "Data truncated for column 'content' at row 1")
  self._do_get_result()
MySQL writing 630 items need 1.8588917255401611 seconds
C:\Users\Windows10\Desktop\code>python es_search.py
ES search data time (seconds): 2.2338857650756836
C:\Users\Windows10\Desktop\code>python mysql_search.py
MySQL search data time (seconds): 0.03143429756164551
C:\Users\Windows10\Desktop\code>
```

结论：Elasticsearch 与 MySQL 比较，在小规模数据情况下，MySQL 的写入与搜索速度更快。