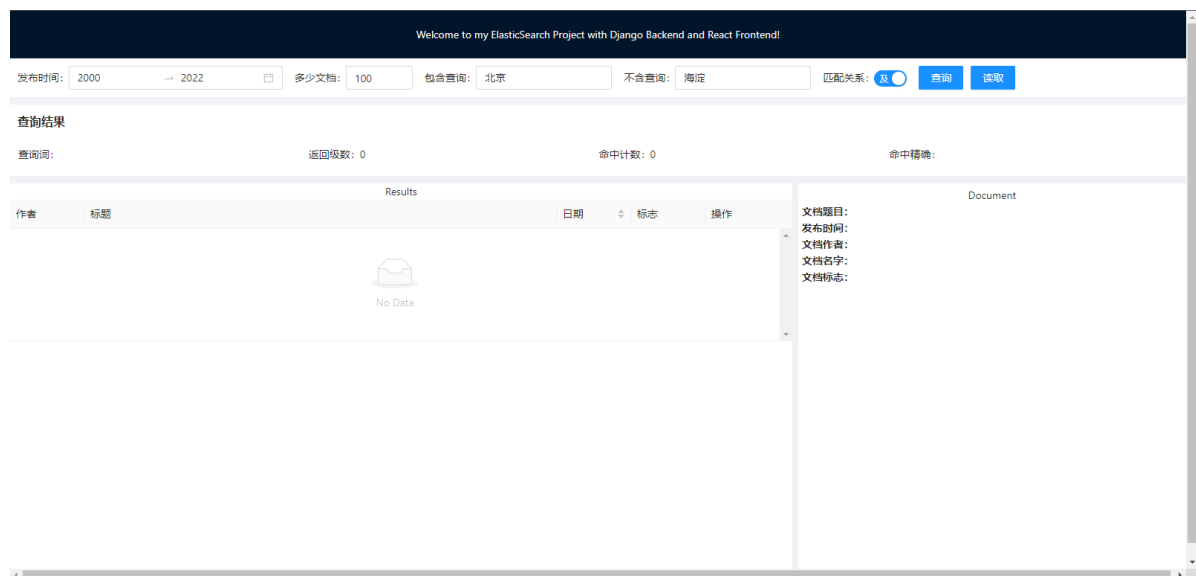# 信息检索 Elasticsearch Pt2

计83 李天勤 2018080106

## Abstract

This homework is an extension of the previous Elasticsearch homework. Our goal add more functionality to our IR system by using more techniques that we had learned in the second half of the semester
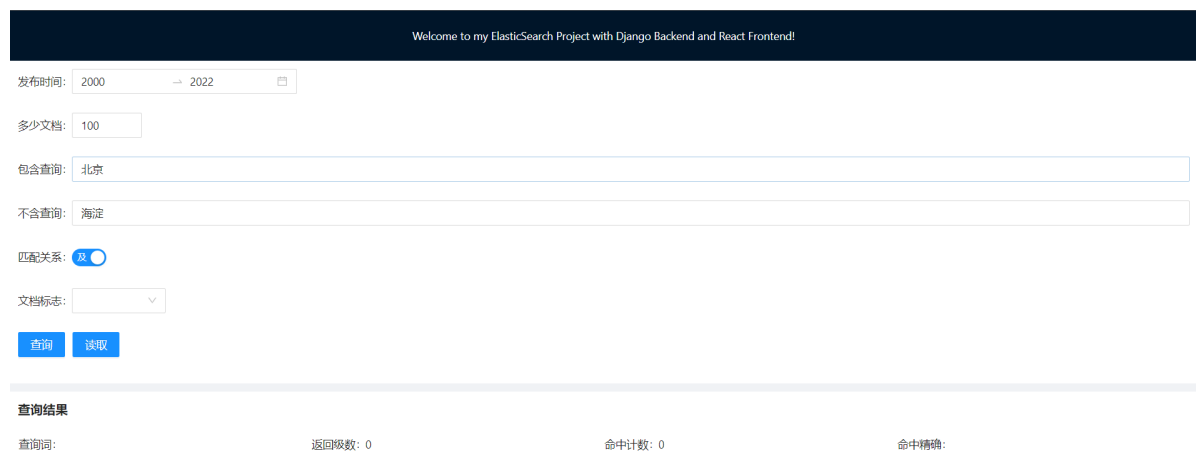
## Background

This is what the experiment initially looked like. It only had basic Elasticsearch functionality which is based on Boolean search. It had two search relationships. If the tab displays 及, then the query has an AND relationship. And the 或 describes an OR relationship. This is explained in the documentation



## Implementation

I slightly changed the format of the page, and added a couple new functionalities. The first small addition is adding functionality that searches documents by its column. It uses Elasticsearch to match the column sections in the index.

匹配关系：及

文档标志：理

理论
理论专页
光辉理论
理论与实践

查询

查询结果

The document table is the same.

| | Results | | | | | |
|---|---|---|---|---|---|---|
| 作者 | 标题 | 评分 | 日期 | 标志 | 操作 | |
| 温红彦;陈晓钟 | 清华大学学生艺术团赴港演出受欢迎 | 8.654198 | 2005-05-07 | 要闻 | 阅读 | |
| | | 8.599146 | 2001-04-24 | | 阅读 | |
| 赵婀娜 | 清华大学将迎来百年校庆 | 8.58141 | 2011-04-09 | 要闻 | 阅读 | |
| 胡和平 | 追求卓越的数学家（纪念华罗庚同志诞辰100周年座谈会发言摘编） | 8.576933 | 2010-11-13 | 综合 | 阅读 | |
| 吴亚明;陈晓星 | 台湾新竹清华大学隆重庆祝清华创校百年 | 8.5559225 | 2011-04-25 | 要闻 | 阅读 | |
| | | 8.555489 | 2001-04-29 | | 阅读 | |
| | | 8.549594 | 2001-04-29 | | 阅读 | |
| 袁新文 | 清华与协和共建医学院 | 8.547539 | 2006-09-06 | 文化新闻 | 阅读 | |
| 赵婀娜 | 水木清华 | 8.534168 | 2011-04-22 | 要闻 | 阅读 | |

Document
文档题目：
发布时间：
文档作者：
文档名字：
文档标志：

< 1 2 > 50 / page

## Word2vec Word Similarity

The following below are the three new functionality that I have added. The first is word similarity. It takes the documents returned by Elasticsearch and does word similarity on the documents.



Word Similarity                                    Word2Vec分析

输入你自己的query或者它会选择elasticsearch返回的第一个文档

Waiting to press button

No Data

Footer

It splits the request in three sections, and returns the time to read, train, and run.

Word Similarity                                    Word2Vec分析

输入你自己的query或者它会选择elasticsearch返回的第一个文档
清华

[0.67025] 北航
[0.62691] 清华学
[0.61717] 北
[0.60058] 北师
[0.58442] 校
[0.58338] 海交
[0.56107] 名校
[0.55162] 南开学
[0.54560] 南开
[0.53846] 高校

○ reading - 39.13464
○ training - 13.21967
○ running - 0.43027

The request on the frontend is as such.

```javascript
trainWordSimilarity = () => {
  this.setState({word_sim_list_items : [], word_sim_pending : "Waiting to read", word_sim_loading : true })
  var result = this.state.documents.map(function(a) {return a.id})
  const data = {"query": this.state.query_include, "document_ids" : result}
  WordSimilarityReading(data).then(response=>{
    const temp = [({"mes" : "reading", "time" : response.data.results})]
    this.setState({ word_sim_reading_time : response.data.results, word_sim_list_items : temp, word_sim_pending : "Waiting to train"})
    console.log(response.data);
  }).then(training => {
    const data2 = {"task" : "train"};
    WordSimilarityTraining(data2).then(response=>{
      const temp = [{"mes": "reading", "time" : this.state.word_sim_reading_time },({"mes" : "training", "time" : response.data.results})]
      this.setState({ word_sim_training_time : response.data.results, word_sim_list_items : temp, word_sim_pending: "Waiting to run"})
      console.log(response.data);
    }).then(running => {
      if (this.state.documents.length === 0) {
        this.setState({ word_sim_pending: "Waiting to press button",word_sim_loading : false, word_sim_list_items : []})
        message.error("没有语料库，请查询")
        return;
      }
      const data3 = {"task" : "run", "query": this.state.word_sim_chosen_word === "" ? this.state.query_include : this.state.word_sim_chosen_word }
      WordSimilarityRunning(data3).then(response=>{
        console.log(response.data);
        if (response.data.results.results === "keyerror") {
          message.error("这个词不在词典里面，重新选")
          this.setState({word_sim_running_time: "", word_sim_list_data: [], word_sim_pending : false, word_sim_loading : false});
          return;
        }
        const temp = [{"mes" : "reading", "time" : this.state.word_sim_reading_time}, {"mes" : "training", "time" : this.state.word_sim_training_time }, ({"mes" : "run
        this.setState({ word_sim_list_items : temp})
        this.setState({word_sim_running_time: response.data.results.time, word_sim_list_data : response.data.results.results, word_sim_pending : false, word_sim_loadi
      }).catch(error=>(message.error(error)))
    }).catch(error=>(message.error(error)))
  }).catch(error=>(message.error(error)))
}
```

And the python backend is simply implements the functionality of word2vec. For example, the train function is as such. The implementations for the following are similar.

```python
def word_sim_train():
    """train corpus based on parameters"""
    start = time.time()
    lines = []
    with open("data.jsonl") as f:
        lines = f.read().splitlines()
    text_corpus = json.loads(lines[0])
    text_corpus = [text.split() for text in text_corpus]

    params = { …
    }

    # train model
    model = Word2Vec(sentences=text_corpus, **params)
    model.save(f"sg{params['sg']}_hs{params['hs']}_win{params['window']}_size{params['vector_size']}.model")

    end = time.time()
    spent_time = end - start
    return '%.5f'%spent_time
```

## Doc2vec Document Similarity

I used Gensim's Doc2vec to implement document similarity. You can enter your own document or query, or it will use the first result returned by elasticsearch. In word2vec, you train to find word vectors and then run similarity queries between words. In doc2vec, you tag your text and you also get tag vectors. Then, after doc2vec training you can use the same vector arithmetic's to run similarity queries on author tags, which document are most similar to the one being queried.

**Document Similarity**                                                    Doc2Vec分析

输入你自己的文档或者它会选择elasticsearch返回的第一个文档

⌐ Waiting to press button

No Data

Footer

It returns the top 5 results

## Document Similarity

输入你自己的文档或者它会选择elasticsearch返回的第一个文档

Doc2Vec分析

[0.32835] [288874] 本报 驻 印尼 记者 董力 一年 前的 5月 27日，印尼 爪哇岛 中部 的 日惹 特别 自治区 发生 了 里氏 5．9 级 地震 和 3 次 里氏 5 级 左右 的 余震，造成 巨大 人员 伤亡 和 财产 损失。在 救灾 和 重建 过程 中，印尼 政府 得到 了 国际 社会 的 大力 援助，印尼 各地 的 华人 社团 也 积极 伸出 援手，受到 当地 政府 和 人民 的 赞许。班 图尔县 是 受灾 最 严重 的 地方，有 4万 多人 死亡，占 死亡 总 人数 的 近 70％；近 14万 幢 房屋 完全 倒塌 或 严重 损毁，占 倒塌 房屋 总数 的 47％。地震 一周年 之际，记者 来到 了 班 图尔县，了解 当地 民众 重建 家园 的 情况。前往 班图尔县 的 路上，还 可以 看到 不少 地震 留下 的 痕迹。不过，在 班图尔县，一幢幢 新 房子 已经 代替 了 去年 的 残垣断壁，很 多 房子 一 看 就是 新修 的。当然，半新半旧 的 房子 也 有 很多，不少 人 仍 在 废墟上 修葺 着 自家 的 房子。大街 上 人来人往，各类 车辆 川流不息。让 人 对 班图尔 两年 重建 计 划 的 按时 完成 充满 信心。班 图尔县 政府 发言人 巴姆邦 介绍 说，班 图尔县 一年 来 的 重建 工作 成果 显著：卫生 医疗 服务 已经 恢复 正常；超过 70％ 的 学校 恢复 教学；道路 和 水利 等 基础 设施 也 基本 完全 修复。巴姆邦 说，班 图尔县 灾后 重建 的 速度 之所以 这么 快，首先 在于 民众 能够 互相 帮助。大家 自发 组织 起来，共同 出力，一家 一家 地 修 复 房屋；其次，离 不开 政府 以及 国际 社会 的 大力 支持。当地 政府 为 每户 灾民 补贴 1500万 印尼盾（相当于 1．5万 元 人民币），联合国 "爪哇 重建 统筹 基金会" 也 为 灾民 提供 大量 资金，用于 房屋 重建，中国 等 国 还 派遣 了 医疗队 救治 在 地震 中 受伤 的 灾民，还有 很多 国家 也 提供 了 援助 物资 和 资金。另外，印尼 各地 的 民间 组织 和 社 团，尤其 是 华人 社团 也 为 重建 做 了 许多 贡献。日惹 福清 公会 主席 陈世详 和 日惹 印尼 华商 总会 秘书 游平茂 介绍 说，日惹 以及 印尼 全国 的 华人 社团 都 纷纷 伸出 援手，参 加 救灾。地震 发生后 不久，日惹 的 华人 社团 就 筹集 了 大批 的 食品、饮用水、毛毯、帐篷 和 药品，并 组织 人员 前往 一个 个 受灾 的 村庄，把 物资 送到 灾民 手上，解 了 他们 的 燃眉之急。日惹 的 华人 社团 还 积极 协助 中国 政府 派来 的 医疗队，为 医疗队 工作 的 顺利 进行 提供 了 便利 条件。灾后 重建 开始后，华人 又 捐赠 了 大量 的 资金，用于 修复 学校、资助 灾民 恢复 生产 等。在 走访 中 记者 看到，日惹 福清 公会 捐资 兴建 的 贾拉坎 小学 校舍 宽敞 明亮（见 上图，本报 记者 董力 摄），是 印尼 华人 社团 捐建 的 近 20 所 学校 中 规模 最大 的 一 所。学校 的 围墙 和 校门 上 写 着 一个 红艳艳 的 "融" 字，校长 穆斯莉曼 女士 介绍 说，这个 "融" 字 代表 了 华人 和 其他 民族 兄弟 般 的 融合 关系。同时，这个 字 也 可以 提醒 学生，长大 以后 要 促进 与 华人 的 团结。陈世详 说，学校 正式 开学 时，日惹特区 领导人 还 在 讲话 中 称赞 华人 社团 的 无私 帮助。班 图尔县 是 整个 日惹 地区 灾后 重建 的 缩影。在 26日 举行 的 地震 一周年 纪念 活动 上，印尼 总统 苏西洛 称赞 日惹 和 中 爪哇 地区 的 所有 组织 工作 都是 由 政府 和 人民 方面 取得 的 成绩，他 不无 骄傲 地 说："以前 外国人 总是 问 印尼 政府 什么 时候 才能 结束 善后 工作，而 现在，他们 问 的 是 政府 如此 迅速 地 进行 灾后 重建 有 什么 秘诀。"苏西洛 的 话 从 另一个 侧面 表明，无论 是 印尼 政府 还是 当地 百姓，都 对 灾区 未来 充满 了 信心。（本报 印尼 日惹电）

# Lsi Document Similarity

Similar to document similarity, but still using LSI.

### Lsi Similarity

输入你自己的文档或者它会选择elasticsearch返回的第一个文档

No Data

Lsi分析

Footer

### Lsi Similarity

输入你自己的文档或者它会选择elasticsearch返回的第一个文档

Lsi分析

[0.9999999] [218980] 本报 香港 5月 6日 电 记者 温红彦、陈晓钟 报道："清韵 华情" 满载 着 促进 清华大学 与 香港 交流 的 使命，满载 着 清华 学子 的 殷殷 祝福，今晚 溢满 了 美丽 的 香江。在 清华大学 校庆 94 周年 之际，清华大学 艺术团 一行 68 位 师生 来到 香港，与 香港 青年 欢度 五四 青年节，为 香港 各界 朋友 和 清华 校友 献上 了 一场场 精彩 纷呈、热 情 奔放 的 民族 歌曲、舞蹈、器乐 和 综合 文艺 节目。清华大学 党委书记 陈希 教授 专程 赴 香港 出席 了 艺术团 的 演出 活动，代表 母校 向 香港 各界 人士、向 广大 的 清华 校友 送 上 节日 问候。清华大学 艺术团 这次 是 应 香港 清华 联会 的 邀请 到 港 演出 的，昨晚 和 今晚 在 香港 理工大学 的 赛马会 综艺 馆 演出，明天 还 将 在 香港 莱涌 袋 锦帆 中学 为 全港 中 学生 演出。"清韵 华情" 是 清华大学 学生 艺术团 的 品牌 型 综合 文艺 演出。中 联办 教育 科技部 副部长 王国力 说，这次 演出 反映 了 清华大学 弘扬 民族 文化、倡导 高雅 艺术 的 宗 旨，展现 了 浓浓 的 中华 情 和 校友 与 母校 之间 的 亲情，同时 也 展示 了 清华 学生 的 艺术 水准。在 港 演出 的 所有 组织 工作 都是 由 香港 清华 联会 的 艺术 教育 成果，以 及 答谢 香港 各界 对 清华大学 建设 的 帮助 和 支持。今晚，清华 学子 们 熟悉 的 蒙民伟 学长 也 观看 了 演出。演出 一 结束，学生们 纷纷 围上 去 与 蒙 先生 合影："天天 进出 清华大学 的 '蒙民伟楼'，今天 终于 见到 了 您！"香港 理工大学 博士生 王树晓、研究生 肖慧 看完 演出 告诉 记者，节目 非常 精彩，反映 了 内地 大学生 昂扬向上 的 精神 风貌，希望 两地 的 大学生 经常 有 这样 的 交流 活动。

[0.830096] [295145] 新华社 记者 刘志杰 作为 庆祝 香港 回归 10 周年 和 纪念 中国 人民 解放军 建军 80 周年 系列 活动 之一，北京 军区 战友 文工团 应 香港 特区 政府 康乐 及 文化事 务署、香港 中华 文化 艺术 基金会 和 香港 合唱团 协会 联合 邀请，8月 25日、26日 与 香港 合唱团 在 港 联合 演出 3 场《长征组歌》大型 合唱 音乐会，演出 获得 极大 成功，演出 前 已 见 多 期待，受 欢迎 演出，是 一次 内地 与 香港、专业 与 业余 相 结合 的 大型 演出 活动。演出 原定 两场，门票 很快 售完，最后 又 加演 一场，门票 同样 销售 一空。演出 前 已 见 多 期待，受 欢迎 盛况。香港 特区 行政长官 曾荫权 观看 了 演出，并 为 音乐会 题词："纪庆 同 欢乐 韵 悠扬"。为 音乐会 提供 独家 赞助 的 香港 中华 文化 艺术 基金会 主席 区永熙 太平绅士 说："能

# Analysis

We can see from the results that Doc2Vec and LSI to compare each method.

Using Elasticsearch, we queried the corpus and got back 10000 results

| 发布时间 | 2000 | → | 2022 | 📅 |
| --- | --- | --- | --- | --- |
| 多少文档 | 10000 | | | |
| 包含查询 | 清华大学 | | | |
| 不含查询 | Ex 海淀 | | | |
| 匹配关系 | 及 🔵 | | | |

The result with the highest score is this Elasticsearch query is

Then, to compare the methods, we can look at the results given. The algorithms and methods used for this extension are the same as the homework for 6 and 7, and the conclusions are pretty much the same. LSI is a count based model where similar terms have same counts for different documents. Then dimensions of this count matrix is reduced using SVD. For both the models similarity can be calculated using cosine similarity. Word2vec is a prediction based model, for example, the given the vector of a word predict the context word vectors (such as skipgram method). When utilizing a small window count, Doc2vec organizes results by terms that are 相似，while LSI organizes results by 相关. Training an LSI system takes much more time than Doc2vec on its basic window of 5, and with 5 epochs.

Compared to Doc2vec and LSI, Elasticsearch is much faster, and cheaper, and handles only pure and simple keyword searches. Thus, when handling simple keyword searches such as 清华 or 北京，Elasticsearch gives us a much better results when we search for documents with simple queries.

Thus, when handling larger document based similarity searches (larger queries), it is better to use Doc2Vec or LSI methods.

# Improvements

Based on this experiment, and throughout this class, it is easy to see how easy to hard it is to make a good IR system. There are a lot of things to consider. If I have more time to work on this experiment, I would like to have implemented more features, and tried more information retrieval methods, such as BERT, and do a deeper comparison between these methods and find a way to display the similarities and differences. I would also like to have improved the interface.

Thanks!