

Evaluation Only. Created with Aspose.Words. Copyright 2003-2016 Aspose Pty Ltd.

LUCENE技文档

说明:本文档在我最大努力范围之内确保其正确性、有效性和可靠性,但并不代表所有的观点都是正确的,而代表个人看法。如有疑问,多指教,谢谢!

Lucene技说明

联系方式:XXXXXXXXXXXX(直接敲本人就好了嘛)

1 技概述

Lucene是apache软件基金会

jakarta项目的一个子项目,是一个开放源代码的全文搜索引擎工具包,但它不是一个完整的全文搜索引擎,而是一个全文搜索引擎的架构,提供了完整的索引引擎和索引引擎,部分文本分析引擎(英文与德文两种西方语言)。

Lucene的目的是为开发人员提供一个易于使用的工具包,以方便的在项目中实现全文搜索的功能,或者是以此为基础建立起完整的全文搜索引擎。Lucene提供了一个但却很大的应用程序接口,能做全文索引和搜索。在Java开发环境里Lucene是一个成熟的免费开源工具。就其本身而言, Lucene是当前以及最近几年最受欢迎的免费Java信息搜索程序。人们常提到信息搜索程序,虽然与搜索引擎有关,但不必将信息搜索程序与搜索引擎相混淆。

Lucene的使用是建立在巨大的索引的基础上的,通常来说,我们在查找时并不直接到数据中去查找,而是直接从索引中查找,因为索引是是二进制格式的,所以搜索速度是很快的。所以我使用Lucene进行索引,要先创建或者取一个索引。

2 使用缺点

2.1 点:

(1)索引文件格式独立于应用平台。Lucene定义了一套以8位字为基础的索引文

件格式,使得兼容系或者不同平台的用能共享建立的索引文件。

(2)在全文搜索引擎的倒排索引的基础上,分了分索引,能新的文件建立小文件索引,提升索引速度。然后通过与原有索引的合并,到的目的。

(3)秀的面向象的系架,使得于Lucene展的学度降低,方便充新功能。

(4)了独立于言和文件格式的文本分析接口,索引器通接受Token流完成索引文件的立,用展新的言和文件格式,只需要文本分析的接口。

(5)已默了一套大的引擎,用无需自己写代即可使系可得大的能力,Lucene的中默了布操作、模糊(Fuzzy Search[11])、分等等。

面已存在的商全文搜索引擎,Lucene也具有相当的。

首先,它的开源代码行方式(遵守Apache Software License[12]),在此基上程序不充分的利用Lucene所提供的大功能,而且可以深入致的学到全文搜索引擎制作技和面向象程的践,而在此基上根据用的情况写出更好的更适合当前用的全文搜索引擎。在一点上,商件的灵活性不及Lucene。

其次,Lucene秉承了开放源代码一的架良的,了一个合理而极具充能力的面向象架,程序可以在Lucene的基上充各种功能,比如充中文理能力,从文本充到HTML、PDF[13]等等文本格式的理,写些展的功能不不复,而且由于Lucene恰当合理的系做了程序上的抽象,展的功能也能易的到跨平台的能力。

最后,移到apache件基金会后,借助于apache件基金会的网平台,程序

可以方便的和开发者、其它程序交流,促成源码的共享,甚至直接得已写完的补充功能。最后,虽然Lucene使用Java写成,但是开放源代码社区的程序员正在不懈的将之使用各种语言(例如.net framework[14]),在遵守Lucene索引文件格式的基础上,使得Lucene能运行在各种各样的平台上,系统管理可以根据当前的平台适合的言来合理的。

2.2缺点

1. Lucene搜索算法不适合网格算

Lucene被写出来的时候硬件没有很大的内存,多处理器也不存在。因此,索引是被做成使用性的内存开很小的方式。即使我花很的来重写跨度算法,并使用多线程内容(使用双核处理器),但是基于迭代器的目标取算法几乎不能。在一些罕的场合你能做一些优化并能迭代一个索引通并行方式,但是大多数场合是不可能的。我遇到的情况是,当我有一个复杂的,超50+的内嵌跨度, CPU在空闲但I/O却一直忙碌,甚至在使用了RAMDirectory.

2. 一个关键的API使得继承Lucene成痛苦

在Lucene的世界中,它被称之为特性。当某些应用需要得到某些,方式是开放的。导致了大多数的都是包保护的,这意味着你不能继承他(除非在你建的似在同一个包下,这样做会污染客代)或者你不得不复制和重写代。更重要的是,如同上面一点提到的,个严重缺乏OO的,一些被内部却没有,匿名被用作复的算当你需要重写他的行。关键API的理由是代在布前得整并且定。虽然想法很光荣,但它再一次让人感到痛苦。因如果你有一些代和Lucene的主要思路并不吻合,你不得不常回Lucene的改到你自己的版本直到

你的补丁被接受。

3. Lucene并非良好设计

Lucene有一个非常糟糕的OO设计。虽然有包, 有接口, 但是它几乎没有任何设计模式。让我想起一个由C(++)开发者的行话, 并且他把坏习惯带到了java中。造成了, 当你需要自定义Lucene来满足你的需求(你将来必定会遇到这样的需求), 你必须面对丑陋的代码。例如:

- 几乎没有使用接口。Query(例如BooleanQuery, SpanQuery, TermQuery...)都是一个抽象类的子类。如果你要添加其中的一个Query, 你会首先想到写一个接口来描述你扩展的契约, 但是抽象的Query并没有定义接口, 你必须经常的实现自己的Query对象到Query中并在本地Lucene中使用。成堆的例子如(HitCollector,...)使用AOP和自代理来也是一个Query。
- 别扭的迭代器。没有hasNext()方法, next()方法返回布尔型并刷新对象内容。你想要保持迭代的元素跟踪来非常的痛苦。我假定它是故意用来节省内存但是它又一次导致了算法上的混乱和复杂。

4. 评分不能被插件化

Lucene有自己评分算法的设计, 当条件增加使用Similarity。但很快它显示出局限性当你想要表示复杂的评分, 例如基于词匹配和元数据的评分。如果你这样做, 你不得不继承Lucene的设计。因为Lucene使用类似tf/idf的评分算法, 然而在我遇到的场合, 在意义上的评分上Lucene的评分机制并不合适。我被迫重写每一个Lucene的评分使得它支持我自定义的评分。它是一个噩梦。

5. 跨度查询太慢

Lingway公司来可能是个特殊的。我的跨度有很要求, Lucene索引已开始添加一, 但它当初可没那么想。最基本的导致了复杂的算法并且行慢, 尤其是当某些短在一份文档中重复了多次出现。是什么我向Lucene是一个高性能的划索引引擎当你使用基本的布。

6. 没有集群的内置支持。

如果你建集群, 你可以写出自己Directory的, 或是使用Solr或者使用Nutch+Hadoop。Solr和Nutch都支持Lucene, 但不是直接的替代。Lucene是可嵌入的, 而你必须支持Solr和Nutch..我Hadoop从Lucene中诞生并不: Lucene并不是通用的。它的内在性决定了大多数合来它是非常快速的, 但是大型文档集合, 你不得不排除Lucene。因它在内核上并没有集群, 你必须把Lucene到搜索引擎, 做并不直接。到Solr或者Nutch上的会你遇到多不必要的麻烦: Nutch中的集成crawling和Solr中的索引。

3 使用步

3.1 前期准

包

Lucene需要的基本包有5个

lucene-core, lucene-analyzers-common, lucene-highlighter, lucene-memory, lucene-queryparse

几个包即可足我在建索引或者搜索的基本需求。

注意: 不同版本的lucene包在操作用的方法也可能会不同, 我下面的例子中使用的版本6.0.0, 其行环境需要1.8以上的JDK。

3.2 建立索引

lucene的“索引”中是存放Document的。Document中存放的是Document的field, 即List<Field>,而我还要用Java的面向对象编程, 不能直接把java的对象直接存放在Lucene的索引中, 所以需要了解下。

建立索引的步骤:

- 1 建立一个IndexWriter对象
- 2 将实体对象转换成Document对象
- 3 将Document对象存入索引中
- 4 关闭源

下面是一个建立索引的例子:

```
package com.lucene.test;

import java.io.File;
import java.io.IOException;
import java.util.ArrayList;
import java.util.List;

import org.apache.lucene.analysis.Analyzer;
import org.apache.lucene.analysis.standard.StandardAnalyzer;
import org.apache.lucene.document.Document;
import org.apache.lucene.document.Field;
import org.apache.lucene.document.FieldType;
import org.apache.lucene.document.TextField;
import org.apache.lucene.index.IndexWriter;
import org.apache.lucene.index.IndexWriterConfig;
import org.apache.lucene.store.Directory;
import org.apache.lucene.store.FSDirectory;
import org.springframework.context.ApplicationContext;
import org.springframework.context.support.ClassPathXmlApplicationContext;

import com.baidu.entity.Item;
import com.baidu.service.ItemService;

public class CreateIndex {
    public static void main(String[] args) {
        ApplicationContext context = new ClassPathXmlApplicationContext("spring.xml");
        ItemService service = context.getBean(ItemService.class);
        //前面部分是查找所有要添加到索引的数据
```

```
List<Item> list = service.findAllItems();

try {
    //索引位置(这个位置可以自定义的)
    Directory directory = FSDirectory.open(new File("E:/dicindex").toPath());
    //分词器
    Analyzer analyzer = new StandardAnalyzer();
    IndexWriterConfig config = new IndexWriterConfig(analyzer);
    //1、创建一个IndexWriter对象
    IndexWriter writer = new IndexWriter(directory,config);

    //创建一个集合存放Document对象
    List<Document> dlist = new ArrayList<Document>();

    //遍历集合list将item挨个存入索引
    for (Item item : list) {
        /*
         * TYPE_STORED      索引,分词,存词
         * TYPE_NOT_STORED 索引,分词,不存词
         */
        Document document = new Document();
        FieldType type = new FieldType(TextField.TYPE_STORED);
        /*type.setStored(false);//是否存词
        type.setTokenized(false);//是否分词*/
        //2、将Article对象转换成document对象
        //第一个参数是搜索词的关键词,可以随意取名
        Field idfield = new Field("id",item.getId().toString(),type);
        Field titlefield = new Field("title",item.getTitle(),type);
        Field contentfield = new Field("content",item.getSellPoint(),type);

        document.add(idfield);
        document.add(titlefield);
        document.add(contentfield);
        dlist.add(document);

    }
    //3、把document对象放入索引中
    writer.addDocuments(dlist);
    //4、关闭源
    writer.close();
} catch (IOException e) {
    e.printStackTrace();
}

System.out.println("索引添加完成");
```

```

    }

}

```

3.3 索引

索引的步骤:

- 1 建立一个IndexSearch对象
- 2 建立一个Query
- 3 将索引出的数据由Document对象体存入集合中
- 4 关闭源

下面是一个从索引中搜索的例子:

```

package com.lucene.test;

import java.io.File;
import java.math.BigInteger;
import java.util.ArrayList;
import java.util.List;
import org.apache.lucene.analysis.Analyzer;
import org.apache.lucene.analysis.standard.StandardAnalyzer;
import org.apache.lucene.document.Document;
import org.apache.lucene.index.DirectoryReader;
import org.apache.lucene.index.IndexReader;
import org.apache.lucene.queryparser.classic.MultiFieldQueryParser;
import org.apache.lucene.queryparser.classic.QueryParser;
import org.apache.lucene.search.IndexSearcher;
import org.apache.lucene.search.Query;
import org.apache.lucene.search.ScoreDoc;
import org.apache.lucene.search.TopDocs;
import org.apache.lucene.store.Directory;
import org.apache.lucene.store.FSDirectory;
import com.baidu.entity.Item;

public class FindIndex {

    public static void main(String[] args) {
        try {
            Directory directory = FSDirectory.open(new File("E:/dicindex").toPath());
            IndexReader reader = DirectoryReader.open(directory);
            IndexSearcher search = new IndexSearcher(reader);
            //分词器
            Analyzer analyzer = new StandardAnalyzer();
            //从n个字段中搜索

```



```
//      QueryParser parser = new QueryParser("content", analyzer);
//      //从多个字段中搜索
//      QueryParser parser = new MultiFieldQueryParser(new String[] { "title", "content" },
analyzer);

Query query = parser.parse("手机");
//第二个参数是指搜索的条数, 我做分词的时候可以把这个数据定大一点, 因为lucene分
词也需要将全部数据都搜索出来再从中找出需要的数据
TopDocs topDocs = search.search(query, 20);
int total = topDocs.totalHits; //搜索出来的文档数
System.out.println("符合关键词(手机)的结果有: "+total);
//搜索出来的文档
ScoreDoc[] scoreDocs = topDocs.scoreDocs;
List<Item> list = new ArrayList<Item>();
for (ScoreDoc scoreDoc : scoreDocs) {
    //关键词的索引
    int index = scoreDoc.doc;
    //根据关键词的索引得到文档的document对象
    Document doc = search.doc(index);
    //将Document对象转换成Item对象
    Item a = new Item();
    a.setId(new BigInteger(doc.get("id")));
    a.setTitle(doc.get("title"));
    a.setSellPoint(doc.get("content"));
    list.add(a);
}
//输出结果
for (Item article : list) {
    System.out.println(article.getTitle());
}
//关键词源
```

This document was truncated here because it was created in the Evaluation Mode.