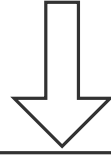
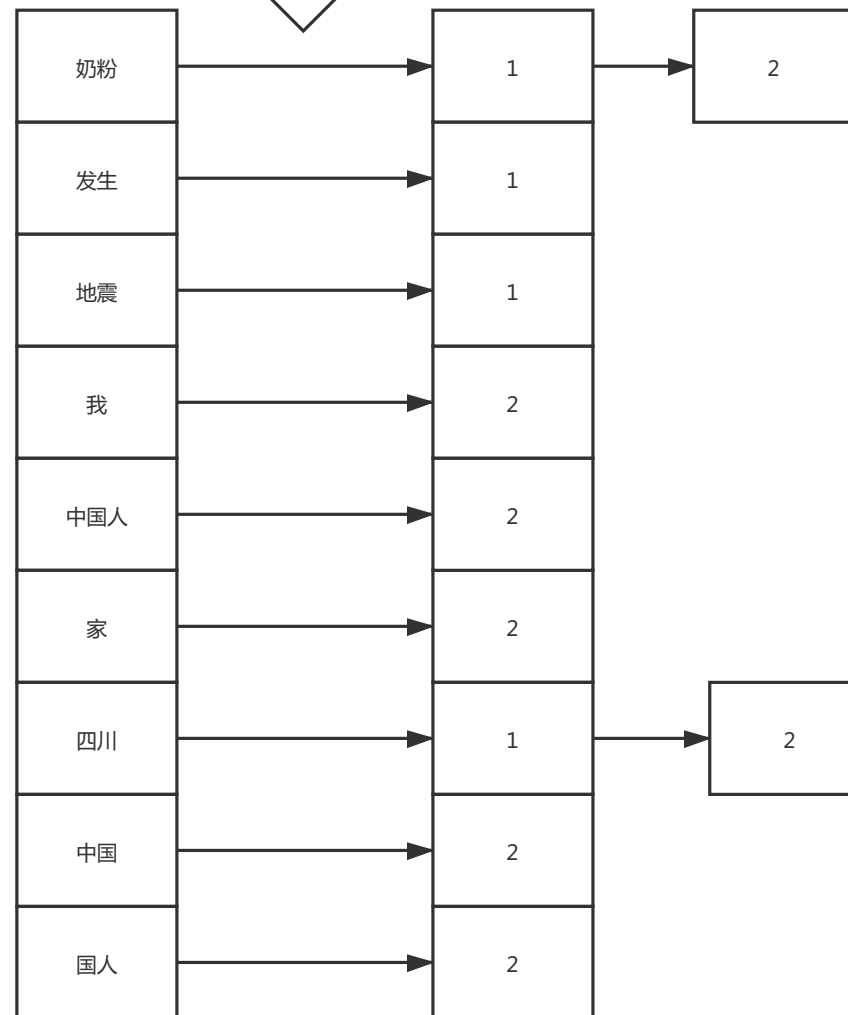
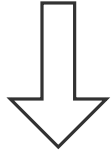


docId : 1, 内容 : 四川省发生了地震  
docId : 2, 内容 : 我是中国人, 我家在四川省



我/中国人/中国/国人/我/家/四川/四川省



my name's zhaoyun ,I like eat  
apples , I am is the xx , running

第一步：空格加标点符号

2.转小写：

3.去停用词

4.会做一个标准化处理

my name zhaoyun I like eat apple  
run

搜apples 会变成搜 apple  
搜apple的时候也是搜 apple

最小颗粒：武汉 武汉市 长江 市长 长江 最长大桥 江大桥 大桥

以上这两个方式怎么用？是不是就是我们都使用最小颗粒去分词呢？那这个smart有什么意义？

smart找不到 我在最小颗粒 那完了 你的倒排索引 得重建。

搜索的时候使用 smart

建索引的用最小颗粒。

为什么都不全采用最小颗粒呢？

为什么不映射：会用在同义词。相同->相似

假设有几千万的文档包含四川。

链表肯定会很长 会产生哪些问题？

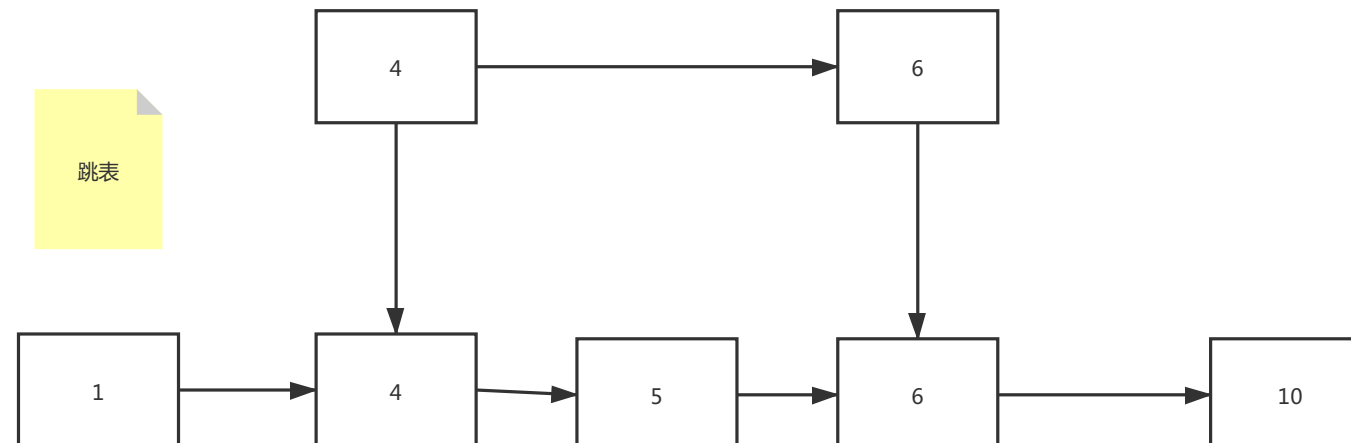
- 1.存储：压缩问题。11122222335444=>3162231534, Lucene的源码解析
- 2.取数据：链表怎么取数据呢？从头遍历 比如有1000个数据 要你拿 第500，遍历500次。怎么办？跳表
- 3.分页问题：

for()

会加载到内存里面来的 让你找第30万。内存会爆炸 这时候通常的做法就是限制100页 或者 根据业务限制的更小。

所有的搜索引擎没有全部数据导出功能

排序：TF-IDF



## 跳表