

# 现代信息检索

## Modern Information Retrieval

第6讲 文档评分、词项权重计算及向量空间模型

Scoring, Term Weighting & Vector Space Model

# 提纲

- 排序式检索
- 词项频率
- tf-idf权重计算
- 向量空间模型

# 布尔检索

- 迄今为止，我们主要关注的是布尔检索，即文档(与查询)要么匹配要么不匹配
- 布尔检索的优点：
  - 对自身需求和文档集性质非常了解的专家而言，布尔查询是不错的选择
  - 对应用开发来说也非常简单，很容易就可以返回1000多条结果
- 布尔检索的不足：
  - 对大多数用户来说不方便
    - 大部分用户不能撰写布尔查询或者他们认为需要大量训练才能撰写出合适的布尔查询
    - 大部分用户不愿意逐条浏览1000多条结果，特别是对Web搜索更是如此

# 布尔检索的其他不足: 结果过少或者过多

- 布尔查询常常会导致过少( $=0$ )或者过多( $>1000$ )的结果
- 例子:
  - 查询 1 (布尔与操作): [standard user dlink 650]
    - $\rightarrow$  200,000 个结果 – 太多
  - 查询2 (布尔与操作): [standard user dlink 650 no card found]
    - $\rightarrow$  0 个结果 – 太少
- 结论:
  - 在布尔检索中, 需要大量技巧来生成一个可以获得合适规模结果的查询

# 排序式检索(Ranked retrieval)

- 排序式检索会对查询和文档的匹配程度进行排序，即给出一个查询和文档匹配评分
- 自由文本查询 (**Free text queries**)：与布尔查询不同，在排序式检索应用中，用户查询通常都是一个或几个关键字
- 排序式检索可以解决返回结果过少或过多的问题
  - 重要的是把相关的结果排在前面

# 排序式检索中的评分技术

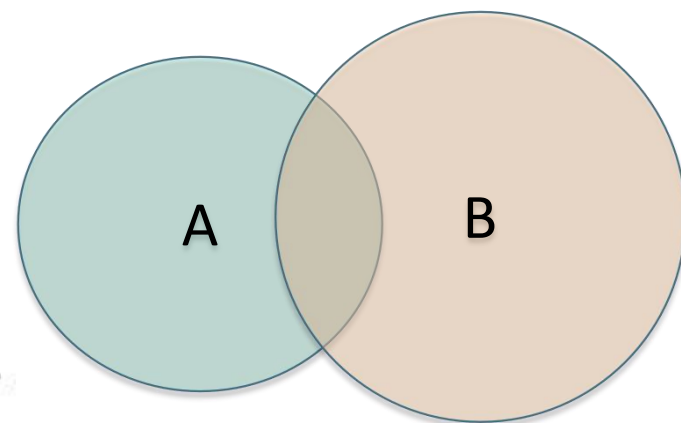
- 我们希望，在同一查询下，文档集中相关度高的文档排名高于相关度低的文档
- 如何实现？
  - 通常做法是对每个查询-文档对赋一个 $[0, 1]$ 之间的分值
  - 该分值度量了文档和查询的匹配程度

# 第一种方法: Jaccard系数

- 计算两个集合重合度的常用方法
- 令  $A$  和  $B$  为两个集合
- Jaccard系数的计算方法:

$$\text{JACCARD}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$(A \neq \emptyset \text{ or } B \neq \emptyset)$$



- $\text{JACCARD}(A, A) = 1$
- $\text{JACCARD}(A, B) = 0$  如果  $A \cap B = \emptyset$
- $A$  和  $B$  不一定要同样大小
- Jaccard 系数会给出一个0到1之间的值

# Jaccard系数的计算样例

- 查询 “ides of March” → 3月15日(朱利乌斯恺撒的殉难日)
- 文档 “Caesar died in March”

$$\text{JACCARD}(q, d) = 1/6$$



# 课堂练习

- 计算下列查询-文档之间的Jaccard系数

q: [information on cars] d1: “all you’ve ever wanted to know about cars”

q: [information on cars] d2: “information on trucks, information on planes, information on trains”

q: [red cars and red trucks] d3: “cops stop red cars more often”

$J(q, d1) = 1/1 = 1$ , cars

$J(q, d2) = 2/6 = 1/3$ , information on

$J(q, d3) = 2/8 = 1/4$ , red cars

可以看出，查询文档词项重合度并不一定与文档相关性一致  
 $J(q, d2)$ 最高但d2与q并不相关：忽略了词项之间重要性的差异。  
Cars比information on更重要

# Jaccard系数的不足

- 不考虑词项频率，即词项在文档中的出现次数(后面会定义)
- 一般而言，罕见词比高频词的信息量更大，Jaccard系数没有考虑这个信息
- 没有仔细考虑文档的长度因素
- 本讲义后面，我们将不考虑使用Jaccard系数来进行查询和文档的评分计算

# 提纲

- 排序式检索
- 词项频率
- tf-idf权重计算
- 向量空间模型

# 查询-文档匹配评分计算

- 如何计算查询-文档的匹配得分？
- 先从单词项查询(查询只包含一个词项)开始
  - 若该词项不出现在文档当中，该文档得分应该为0
  - 该词项在文档中出现越多，则得分越高
  - 这就是所谓**词项频率** (term frequency, 简称tf)评分
- 后面我们将给出多种评分的方法

# 二值关联矩阵

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbet h ...
--	-----------------------------	------------------	----------------	--------	---------	-----------------

ANTHONY	1	1	0	0	0	1
BRUTUS	1	1	0	1	0	0
CAESAR	1	1	0	1	1	1
CALPURNIA	0	1	0	0	0	0
CLEOPATRA	1	0	0	0	0	0
MERCY	1	0	1	1	1	1
WORSER	1	0	1	1	1	0
...						

每篇文档可以看成是一个二值的向量  $\in \{0, 1\}^{|V|}$

# 非二值关联矩阵(词项频率)

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth ...
ANTHONY	157	73	0	0	0	1
BRUTUS	4	157	0	2	0	0
CAESAR	232	227	0	2	1	0
CALPURNIA	0	10	0	0	0	0
CLEOPATRA	57	0	0	0	0	0
MERCY	2	0	3	8	5	8
WORSER	2	0	1	1	1	5
...						

每篇文档可以表示成一个词频向量  $\in \mathbb{N}^{|V|}$

# 词袋(Bag of words)模型

- 不考虑词在文档中出现的顺序
- *John is quicker than Mary* 及 *Mary is quicker than John* 的表示结果一样
- 这称为一个词袋模型(bag of words model, BOW模型)
- 在某种意义上说，这种表示方法是一种“倒退”，因为位置索引中能够区分上述两篇文档

# 词项频率 tf

- 词项 $t$ 的词项频率(以下简称词频)  $tf_{t,d}$  是指 $t$  在 $d$ 中出现的次数, 是与文档相关的一个量, 可以认为是**文档内代表度**的一个量, 也可以认为是一种**局部信息**。
- 下面将介绍利用 $tf$ 来计算文档评分的方法
- 第一种方法是采用原始的 $tf$ 值(raw  $tf$ )
- 但是原始 $tf$ 不太合适:
  - 某个词项在A文档中出现十次, 即 $tf = 10$ , 在B文档中 $tf = 1$ , 那么A比B更相关
  - 但是相关度不会相差10倍, 即相关度不会正比于词项频率 $tf$



# 一种替代原始tf的方法: 对数词频

- $t$  在  $d$  中的对数词频权重定义如下:

$$w_{t,d} = \begin{cases} 1 + \log_{10} \text{tf}_{t,d} & \text{if } \text{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$

- $\text{tf}_{t,d} \rightarrow w_{t,d}$ :  
 $0 \rightarrow 0, 1 \rightarrow 1, 2 \rightarrow 1.3, 10 \rightarrow 2, 1000 \rightarrow 4$ , 等等
- 文档-词项的匹配得分是所有同时出现在 $q$ 和文档 $d$ 中的词项的对数词频之和  $\sum_{t \in q \cap d} (1 + \log \text{tf}_{t,d})$
- 如果两者没有公共词项, 则得分为0

# 提纲

- 排序式检索
- 词项频率
- tf-idf权重计算
- 向量空间模型

# 文档集中的词频 vs. 文档中的词频

单词	文档集频率	文档频率
INSURANCE	10440	3997
TRY	10422	8760

- 词项 $t$ 的文档集频率(Collection frequency,  $cf$ ): 文档集中出现的 $t$ 词条的个数
- 词项 $t$ 的文档频率 $df$ : 包含 $t$ 的文档篇数
- 为什么会出现上述表格的情况? 即文档集频率相差不大, 但是文档频率相差很大
- 哪个词是更好的搜索词项? 即应该赋予更高的权重
- 上例表明  $df$  (和 $idf$ ) 比 $cf$  (和“ $icf$ ”)更适合权重计算

# 罕见词项所期望的权重

- 罕见词项比常见词所蕴含的信息更多
- 考虑查询中某个词项，它在整个文档集中非常罕见 (例如 ARACHNOCENTRIC).
- 某篇包含该词项的文档很可能相关
- 于是，我们希望像ARACHNOCENTRIC一样的罕见词项将有较高权重
- 物以稀为贵！

# 常见词项所期望的权重

- 常见词项的信息量不如罕见词
- 考虑一个查询词项，它频繁出现在文档集中 (如 GOOD, INCREASE, LINE等等)
- 一篇包含该词项的文档当然比不包含该词项的文档的相关度要高
- 但是，这些词对于相关度而言并不是非常强的指示词
- 于是，对于诸如GOOD、INCREASE和LINE的频繁词，会给一个正的权重，但是这个权重小于罕见词权重

# 文档频率(Document frequency, df)

- 对于罕见词项我们希望赋予高权重
- 对于常见词我们希望赋予正的低权重
- 接下来我们使用文档频率df这个因子来计算查询-文档的匹配得分
- 文档频率(document frequency, df)指的是出现词项的文档数目

# idf 权重

- $df_t$  是出现词项 $t$ 的文档数目
- $df_t$  是和词项 $t$ 的信息量成反比的一个值
- 于是可以定义词项 $t$ 的idf权重(逆文档频率):

$$idf_t = \log_{10} \frac{N}{df_t}$$

(其中 $N$  是文档集中文档的数目)

- $idf_t$  是反映词项 $t$ 的信息量的一个指标，是一种全局性指标，反应的是词项在全局的区别性。
- 实际中往往计算 $[\log N/df_t]$ 而不是  $[N/df_t]$ ，这可以对idf的影响有所抑制
- 值得注意的是，对于tf 和idf我们都采用了对数计算方式

# idf 的另一种解读

$$\text{Log } N/df = -\log df/N = -\log P(t \mid \text{Collection})$$

$-\log P(E)$ ：事件E的信息量

因而，idf可视为一种词项全局信息量的指标



# idf的计算样例

- 利用右式计算idf<sub>t</sub>:

$$\text{idf}_t = \log_{10} \frac{1,000,000}{\text{df}_t}$$

词项	df <sub>t</sub>	idf <sub>t</sub>
calpurnia	1	6
animal	100	4
sunday	1000	3
fly	10,000	2
under	100,000	1
the	1,000,000	0

假设语料中文档数目N=1, 000, 000

# idf对排序的影响

- 对于单词项查询,idf对文档排序没有任何影响
- idf 会影响至少包含2个词项的查询的文档排序结果
- 例如, 在查询 “arachnocentric line” 中, idf权重计算方法会增加arachnocentric的相对权重, 同时降低 line的相对权重

# tf-idf权重计算

- 词项的tf-idf权重是tf权重和idf权重的乘积

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \cdot \log \frac{N}{\text{df}_t}$$

- 信息检索中最出名的权重计算方法
- 其他叫法：tf.idf、tf x idf

# tf-idf小结

- 词项t在文档d中的权重可以采用下式计算

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \cdot \log \frac{N}{\text{df}_t}$$

- tf-idf权重
  - 随着词项频率的增大而增大 (局部信息)
  - 随着词项罕见度的增加而增大 (全局信息)

# 课堂练习: 词项、文档集及文档频率

统计量	符号	定义
词项频率	$tf_{t,d}$	$t$ 在文档 $d$ 中出现的次数
文档频率	$df_t$	出现 $t$ 的文档数目
文档集频率	$cf_t$	$t$ 在文档集中出现的总次数

- $df$ 和 $cf$ 有什么关系?
- $tf$ 和 $cf$ 有什么关系?
- $tf$ 和 $df$ 有什么关系?

# 提纲

- 排序式检索
- 词项频率
- tf-idf权重计算
- 向量空间模型

# 二值关联矩阵

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth ...
ANTHONY	1	1	0	0	0	1
BRUTUS	1	1	0	1	0	0
CAESAR	1	1	0	1	1	1
CALPURNIA	0	1	0	0	0	0
CLEOPATRA	1	0	0	0	0	0
MERCY	1	0	1	1	1	1
WORSER	1	0	1	1	1	0
...						

每篇文档表示成一个二值向量  $\in \{0, 1\}^{|V|}$

# tf矩阵

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbet h ...
ANTHONY	157	73	0	0	0	1
BRUTUS	4	157	0	2	0	0
CAESAR	232	227	0	2	1	0
CALPURNIA	0	10	0	0	0	0
CLEOPATRA	57	0	0	0	0	0
MERCY	2	0	3	8	5	8
WORSE	2	0	1	1	1	5
...						

每篇文档表示成一个词频向量  $\in \mathbb{N}^{|V|}$



# 二值 $\rightarrow$ tf矩阵 $\rightarrow$ tfidf矩阵

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbet h ...
--	-----------------------------	------------------	----------------	--------	---------	-----------------

ANTHONY	5.25	3.18	0.0	0.0	0.0	0.35
BRUTUS	1.21	6.10	0.0	1.0	0.0	0.0
CAESAR	8.59	2.54	0.0	1.51	0.25	0.0
CALPURNIA	0.0	1.54	0.0	0.0	0.0	0.0
CLEOPATRA	2.85	0.0	0.0	0.0	0.0	0.0
MERCY	1.51	0.0	1.90	0.12	5.25	0.88
WORSER	1.37	0.0	0.11	4.15	0.25	1.95
...						

每篇文档表示成一个基于tfidf权重的实值向量  $\in \mathbb{R}^{|V|}$

# 文档表示成向量

- 每篇文档表示成一个基于tfidf权重的实值向量  $\in \mathbb{R}^{|V|}$ .
- 于是，我们有一个  $|V|$ 维实值空间
- 空间的每一维都对应词项
- 文档都是该空间下的一个点或者向量
- 极高维向量：对于Web搜索引擎，空间会上千万维
- 对每个向量来说又非常稀疏，大部分都是0

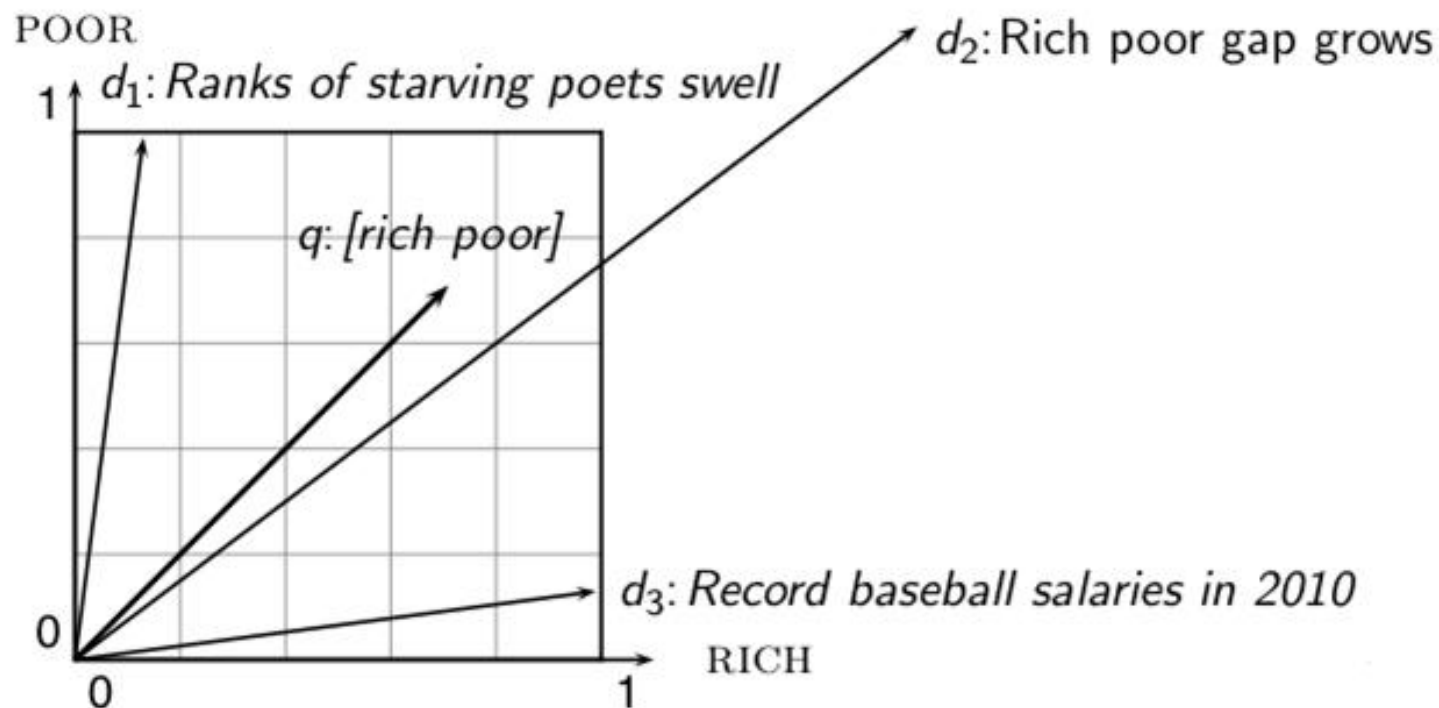
# 查询表示成向量

- 关键思路1: 对于查询做同样的处理, 即将查询表示成同一高维空间的向量
- 关键思路2: 按照文档对查询的邻近程度排序
  - 邻近度 = 相似度
  - 邻近度  $\approx$  距离的反面
- 回想一下, 我们是希望和布尔模型不同, 能够得到非二值的、既不是过多或也不是过少的检索结果
- 这里, 我们希望相关文档的相关度评分高于不相关文档

# 向量空间下相似度的形式化定义

- 先考虑一下两个点之间的距离的倒数
- 一种方法是采用欧氏距离
- 但是，欧氏距离不是一种好的选择，这是因为欧氏距离对向量长度很敏感

# 欧氏距离不好的例子



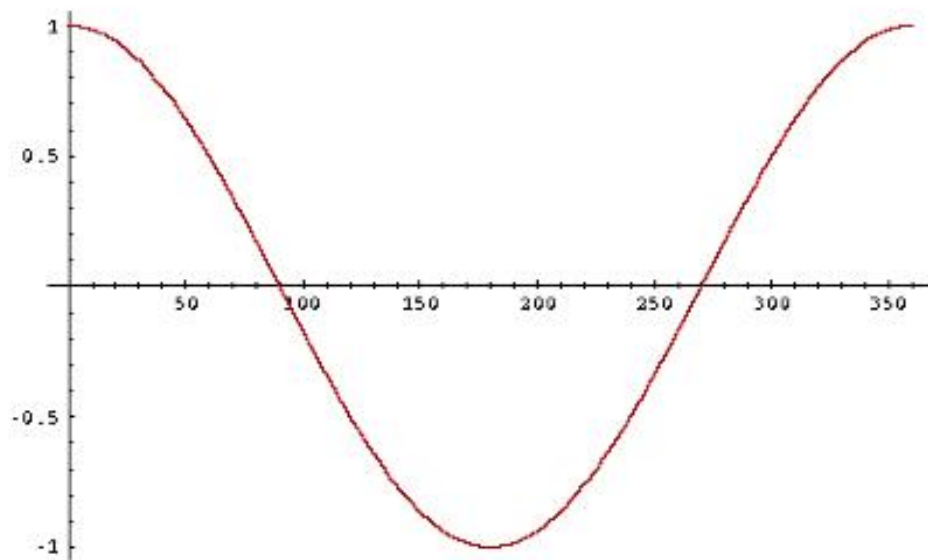
尽管查询 $q$ 和文档 $d_2$ 的词项分布非常相似，但是采用欧氏距离计算它们对应向量之间的距离非常大。

# 采用夹角而不是距离来计算

- 将文档按照其向量和查询向量的夹角大小来排序
- 假想实验：将文档  $d$  复制一份加在自身末尾得到文档  $d'$ .  
 $d'$  是  $d$  的两倍
- 很显然，从语义上看， $d$  和  $d'$  具有相同的内容
- 两者之间的夹角为0，代表它们之间具有最大的相似度
- 但是，它们的欧氏距离可能会很大

# 从夹角到余弦

- 下面两个说法是等价的：
  - 按照夹角从小到大排列文档
  - 按照余弦从大到小排列文档
- 这是因为在区间 $[0^\circ, 180^\circ]$ 上，余弦函数cosine是一个单调递减函数



# 查询和文档之间的余弦相似度计算

$$\cos(\vec{q}, \vec{d}) = \text{SIM}(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}| |\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

- $q_i$  是第  $i$  个词项在查询  $q$  中的 tf-idf 权重
- $d_i$  是第  $i$  个词项在文档  $d$  中的 tf-idf 权重
- $|\vec{q}|$  和  $|\vec{d}|$  分别是  $\vec{q}$  和  $\vec{d}$  的长度
- 上述公式就是  $\vec{q}$  和  $\vec{d}$  的余弦相似度，或者说向量  $\vec{q}$  和  $\vec{d}$  的夹角的余弦



# 文档长度归一化

- 如何计算余弦相似度？
- 一个向量可以通过除以它的长度进行归一化处理，以下使用 $L_2$ （2范数）：

$$\|x\|_2 = \sqrt{\sum_i x_i^2}$$

- 这相当于将向量映射到单位球面上
- 这是因为归一化之后： $\|x\|_2 = \sqrt{\sum_i x_i^2} = 1.0$
- 因此，长文档和短文档的向量中的权重都处于同一数量级
- 前面提到的文档  $d$  和  $d'$  (两个  $d$  的叠加) 经过上述归一化之后的向量相同

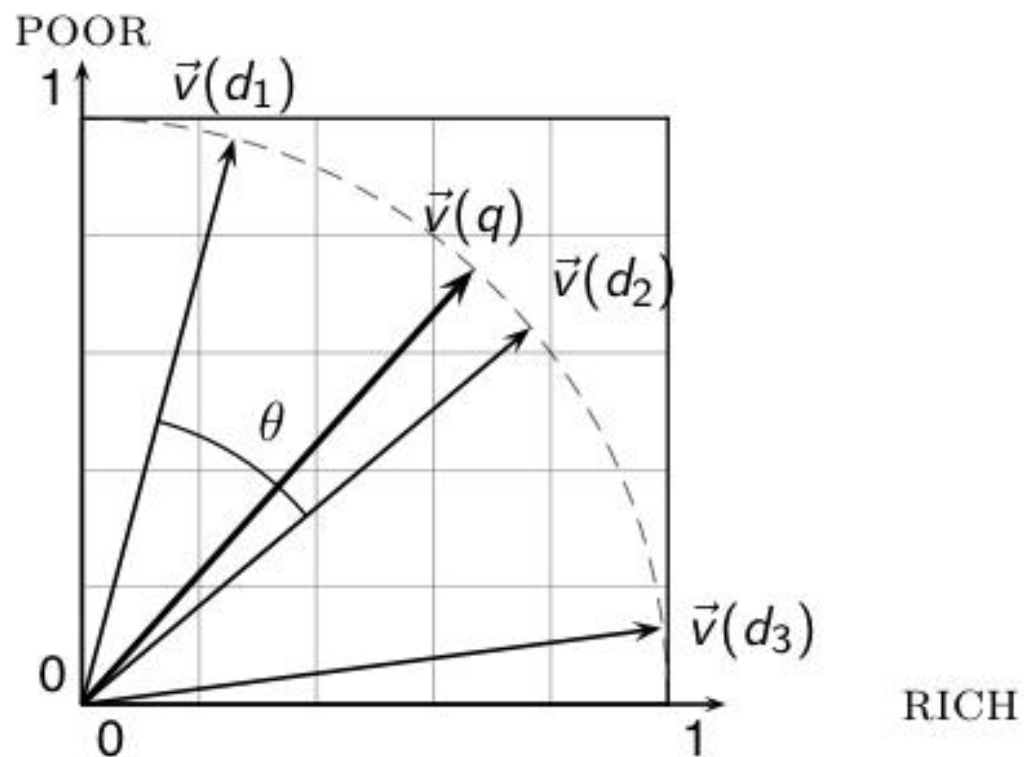
# 归一化向量的余弦相似度

- 归一化向量的余弦相似度等价于它们的点积(或内积)

$$\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d} = \sum_i q_i \cdot d_i$$

如果  $\vec{q}$  和  $\vec{d}$  都是长度归一化后的向量

# 余弦相似度的图示



# 余弦相似度的计算样例

词项频率tf

3本小说之间的相似度

- (1) SaS(理智与情感):  
Sense and  
Sensibility
- (2) PaP(傲慢与偏见):  
Pride and  
Prejudice
- (3) WH(呼啸山庄):  
Wuthering  
Heights

词项	SaS	PaP	WH
AFFECTION	115	58	20
JEALOUS	10	7	11
GOSSIP	2	0	6
WUTHERING	0	0	38

# 余弦相似度计算

词项频率 tf

词项	SaS	PaP	WH
AFFECTION	115	58	20
JEALOUS	10	7	11
GOSSIP	2	0	6
WUTHERING	0	0	38

对数词频 ( $1+\log_{10}tf$ )

词项	SaS	PaP	WH
AFFECTION	3.06	2.76	2.30
JEALOUS	2.0	1.85	2.04
GOSSIP	1.30	0	1.78
WUTHERING	0	0	2.58

为了简化计算，上述计算过程中没有引入idf

# 余弦相似度计算

对数词频 ( $1 + \log_{10} \text{tf}$ )

词项	SaS	PaP	WH
AFFECTION	3.06	2.76	2.30
JEALOUS	2.0	1.85	2.04
GOSSIP	1.30	0	1.78
WUTHERING	0	0	2.58

对数词频的余弦归一化结果

词项	SaS	PaP	WH
AFFECTION	0.789	0.832	0.524
JEALOUS	0.515	0.555	0.465
GOSSIP	0.335	0.0	0.405
WUTHERING	0.0	0.0	0.588

$$\cos(\text{SaS}, \text{PaP}) \approx 0.789 * 0.832 + 0.515 * 0.555 + 0.335 * 0.0 + 0.0 * 0.0 \approx 0.94.$$

$$\cos(\text{SaS}, \text{WH}) \approx 0.79$$

$$\cos(\text{PaP}, \text{WH}) \approx 0.69$$

$$\cos(\text{SaS}, \text{PaP}) > \cos(\text{SaS}, \text{WH}) > \cos(\text{PaP}, \text{WH})$$

# 余弦相似度计算算法

COSINESCORE( $q$ )

```
1  float Scores[ $N$ ] = 0
2  float Length[ $N$ ]
3  for each query term  $t$ 
4  do calculate  $w_{t,q}$  and fetch postings list for  $t$ 
5      for each pair( $d, tf_{t,d}$ ) in postings list
6      do  $Scores[d] + = w_{t,d} \times w_{t,q}$ 
7  Read the array Length
8  for each  $d$ 
9  do  $Scores[d] = Scores[d] / Length[d]$ 
10 return Top  $K$  components of Scores[]
```

假设数组Length预先存放了文档  
归一化因子，即二范式

# 长度归一化

---

- 考虑三个因素tf、idf、文档长度
- 为什么长度因素很重要？
  - 长度长的文档词项也多
  - 长度长的文档TF高
- 需要对长度进行规整/归一(normalization)处理！



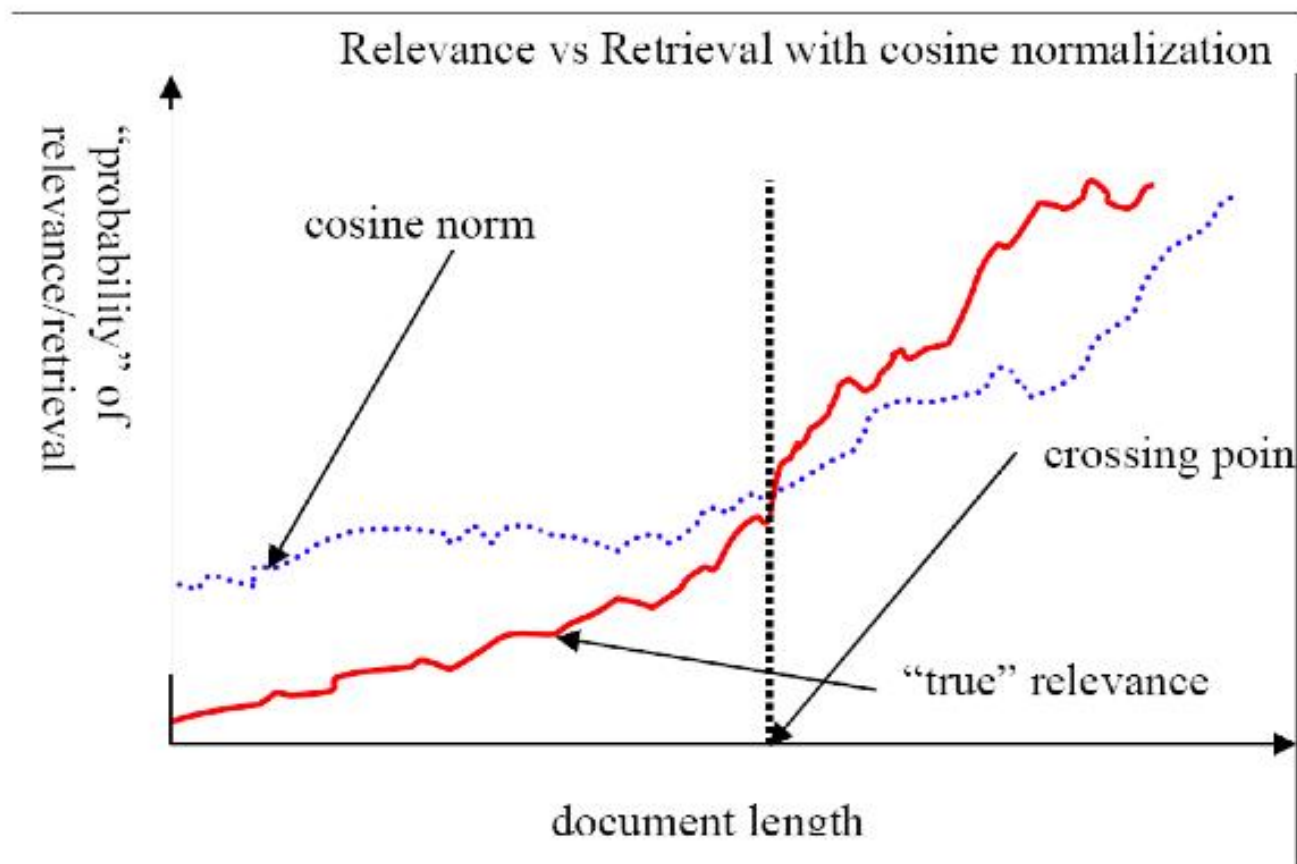
# 课堂练习：余弦相似度的一个问题

- 查询 q: “anti-doping rules Beijing 2008 olympics”
- 计算并比较如下的三篇文档
  - d1: 一篇有关 “anti-doping rules at 2008 Olympics” 的短文档
  - d2: 一篇包含 d1 以及其他 5 篇新闻报道的长文档，其中这 5 篇新闻报道的主题都与 Olympics/anti-doping 无关
  - d3: 一篇有关 “anti-doping rules at the 2004 Athens Olympics” 的短文档
- 我们期望的结果是什么？  $d1 > d2 > d3$
- 但是实际上按余弦相似度排序是  $d1 > d3 > d2$
- 余弦相似度可能对长度过度惩罚

# 回转归一化

- 余弦归一化倾向于短文档，即对短文档产生的归一化因子太大，而平均而言对长文档产生的归一化因子太小
  - 注：这里“归一化因子”指的是2范数的倒数，即将“因子”视为一个系数而不是分母
  - 余弦归一化对长文档的惩罚过重，实际上长文档中虽然词频较高，但也会包含较多的信息
- 于是可以先找到一个支点(pivot, 平衡点)，然后通过这个支点对余弦归一化操作进行线性调整。
- 效果：短文档的相似度降低，而长文档的相似度增大
- 这可以去除原来余弦归一化偏向短文档的问题

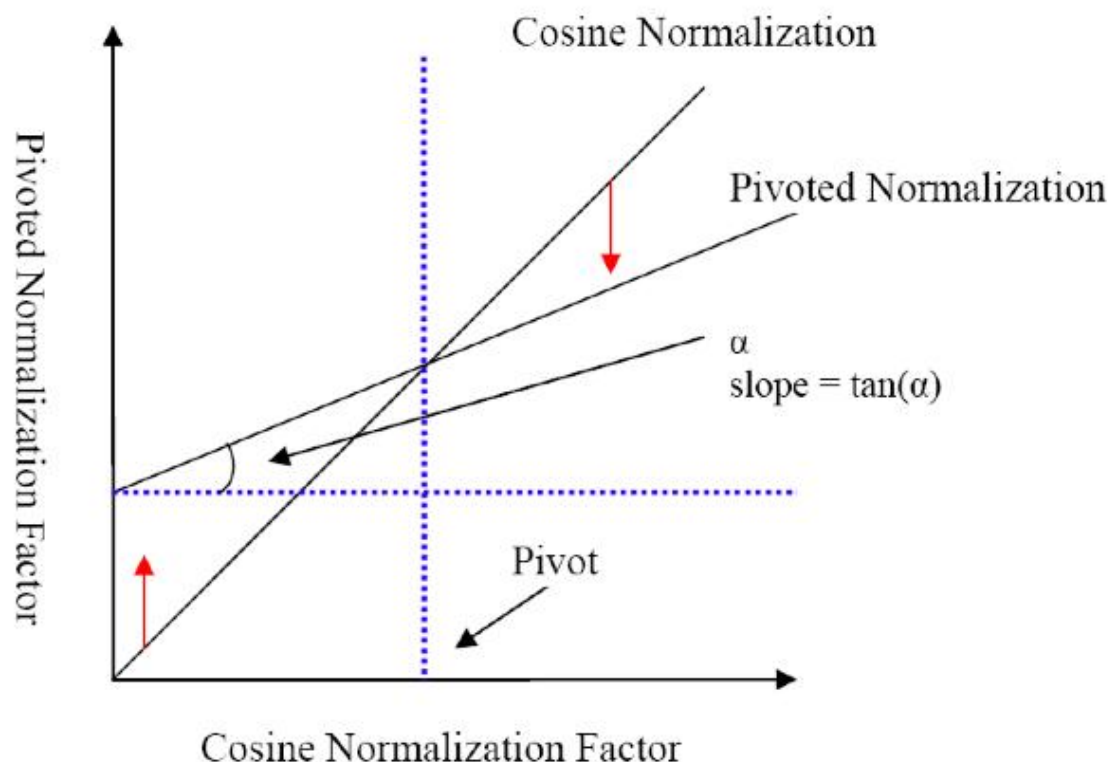
# 预测相关性概率 vs. 真实相关性概率



- X轴是文档长度，Y轴是文档“相关”的概率。这个概率从IR评测语料计算。
- 红色曲线是实际概率分布，蓝色虚线是由cosine归一化估计的概率分布。可见两条曲线存在一定的差异。

# 回转归一化(Pivot normalization)

Pivot normalization



- 基本思想：旋转归一化曲线，使得两条曲线尽量重合

# 回转归一化: Amit Singhal的实验结果

Cosine	Pivoted Cosine Normalization				
	Slope				
	0.60	0.65	0.70	<b>0.75</b>	0.80
6,526	6,342	6,458	6,574	<b>6,629</b>	6,671
0.2840	0.3024	0.3097	0.3144	<b>0.3171</b>	0.3162
Improvement	+ 6.5%	+ 9.0%	+10.7%	<b>+11.7%</b>	+11.3%

- 结果第一行：返回的相关文档数目
- 结果第二行： 平均正确率
- 结果第三行： 平均正确率的提高百分比

Slope: 斜率, 即“回转”参数  
baseline是Cosine归一化

# Amit Singhal

- 1989年本科毕业于印度IIT (Indian Institute of Technology) Roorkee分校
- 1996年博士毕业于Cornell University，导师是Gerard Salton
- 其论文获得1996年SIGIR Best Student Paper Award

[Pivoted Document Length Normalization](http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.50.9950) - [ 翻译此页 ]  
citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.50.9950 - 网页快照

作者: A Singhal - 1996 - 被引用次数: 716 - 相关文章

Document length normalization is used to fairly retrieve documents ...

- 2551 – Introduction to Modern Information retrieval – Salton, McGill - 1983
- 1078 – Term weighting approaches in automatic text retrieval – Salton, Buckley ...
- 203 – Inference networks for document retrieval – Turtle, Croft - 1990

## The smartest people in tech

Engineer runner-up: Amit Singhal

24 of 50 Back Next

Google Fellow,  
Engineering

If you've ever used Google search, you've benefited from Singhal's work. The software engineer heads up the company's core ranking team, which is constantly working to improve the accuracy, speed, and thoroughness of Google search. When Singhal started at Google in 2000, the search engine of today was pure science fiction. The time between the posting of a web page to the time a person could find that page via Google could take anywhere from 15 to 45 days. Now, thanks to Singhal and his



• How we chose the smartest people in tech

- 2000年加入Google，2001年被授予Google Fellow称号
- Google 排序团队负责人,被财富杂志(Fortune, 2010)誉为世界科技界最聪明的50个人之一

# tf-idf权重计算的三要素

词项频率tf		文档频率df		归一化方法	
n(natural)	$tf_{t,d}$	n(no)	1	n(none)	1
l(logarithm)	$1 + \log(tf_{t,d})$	t(idf)	$\log \frac{N}{df_t}$	c(cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a(augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_t (tf_{t,d})}$	p(prob idf)	$\max \left\{ 0, \log \frac{N - df_t}{df_t} \right\}$	u(pivoted unique)	$1/u(17.4.4节)$
b(boolean)	$\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			b(byte size)	$1/CharLength^a, a < 1$
L(log ave)	$\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$				

权重机制：q和d的tf-idf权重是三个因子的乘积，每个因子的计算方法有多种选项，如上表。q和d的权重计算方式可以不同（例如，q可以不采用归一化，即归一化方法n）。最终的文档评分是归一化后的q、d向量的内积

# tf-idf权重机制举例

- 对于查询和文档常常采用不同的权重计算机制
- 记法: ddd.qqq
- 例如: Inc.ltn
  - 文档: 对数tf, 无idf因子, 余弦长度归一化
  - 查询: 对数tf, idf, 无归一化
- 文档当中不用idf结果会不会很差?
  - 查询: “best car insurance”
  - 文档: “car insurance auto insurance”
  - 一般来说不会, 因为查询中使用了idf因子。最终文档得分是d、q权重的内积, 已经乘了idf。



# tf-idf 计算样例: Inc.ltn

查询: “best car insurance”. 文档: “car insurance auto insurance”.

word	query					document				product
	tf-raw	tf-wght	df	idf	weight	tf-raw	tf-wght	weight	n'lized	
auto	0	0	5000	2.3	0	1	1	1	0.52	0
best	1	1	50000	1.3	1.3	0	0	0	0	0
car	1	1	10000	2.0	2.0	1	1	1	0.52	1.04
insurance	1	1	1000	3.0	3.0	2	1.3	1.3	0.68	2.04

$$\sqrt{1^2 + 0^2 + 1^2 + 1.3^2} \approx 1.92$$

$$1/1.92 \approx 0.52$$

$$1.3/1.92 \approx 0.68$$

$$\text{最终结果 } \sum w_{qi} \cdot w_{di} = 0 + 0 + 1.04 + 2.04 = 3.08$$

# 向量空间模型小结

- 将查询表示成tf-idf权重向量
  - 将每篇文档表示成同一空间下的 tf-idf权重向量
  - 计算两个向量之间的某种相似度(如余弦相似度)
  - 按照相似度大小将文档排序
  - 将前 $K$ （如 $K=10$ ）篇文档返回给用户
- 
- 向量空间(检索)模型中包含一个向量表示模型，可以广泛用于其他领域，比如：判断论文抄袭、代码抄袭、图像相似度计算等等

# Gerard Salton(1927-1995)

- 信息检索领域的奠基人之一，向量空间模型的完善者和倡导者，SMART系统的主要研制者，ACM Fellow
- 1958年毕业于哈佛大学应用数学专业，是Howard Aiken的关门博士生。Howard Aiken是IBM第一台大型机ASCC的研制负责人。
- Salton是康奈尔大学计算机系的创建者之一。



# 本讲内容

- 对搜索结果排序(Ranking) : 为什么排序相当重要?
- 词项频率(Term Frequency, TF): 排序中的重要因子
- Tf-idf 权重计算方法: 最出名的经典排序方法
- 向量空间模型(Vector space model): 信息检索中最重要的形式化模型之一 (其他模型还包括布尔模型和概率模型)

# 本讲小结

- tf, 反映词项在文档内部的权重, 局部权重
- idf, 反映词项在文档集(文档间)的权重, 全局权重
- tf-idf, 综合以后反映词在文档中的权重
- 长度因素, 归一化

# 参考资料

- 《信息检索导论》第6、7章
- <http://ifnlp.org/ir>
  - 向量空间入门
  - Exploring the similarity space (Moffat and Zobel, 2005)
  - Okapi BM25 (另外一种著名的权重计算方法, 《信息检索导论》11.4.3节)

# 课后练习

---

- 有待补充