
NLP文本相似度

Outline

余弦相似度、向量空间模型

TFIDF

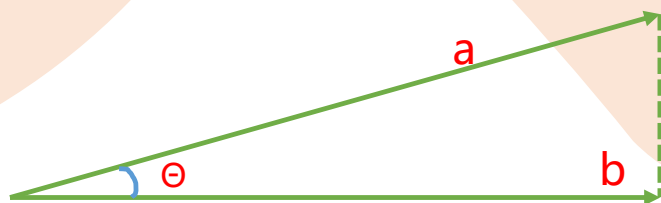
【实践】TFIDF

LCS

【实践】LCS

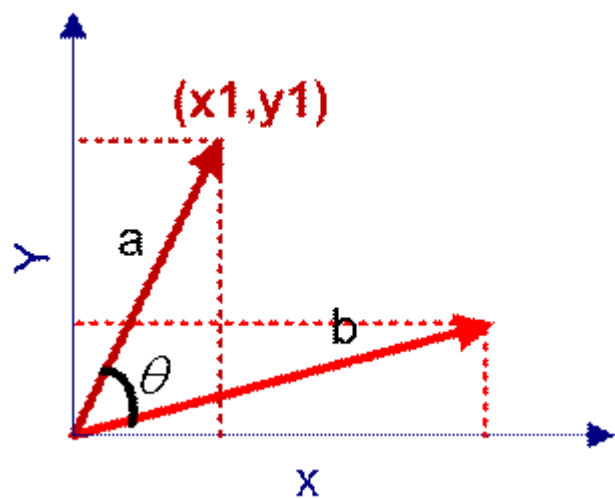
相似度

- 相似度度量：计算个体间相似程度
- 相似度值越小，距离越大，相似度值越大，距离越小
- 最常用——余弦相似度
 - 一个向量空间中两个向量夹角的余弦值作为衡量两个个体之间差异的大小
 - 余弦值接近1，夹角趋于0，表明两个向量越相似



$$\cos(\theta) = \frac{b}{a}$$

相似度



如果向量a和b不是二维而是n维

$$\cos(\theta) = \frac{\sum_{i=1}^n (x_i \times y_i)}{\sqrt{\sum_{i=1}^n (x_i)^2} \times \sqrt{\sum_{i=1}^n (y_i)^2}}$$
$$= \frac{a \bullet b}{||a|| \times ||b||}$$

一个例子

- 句子1：这只皮靴号码大了，那只号码合适
- 句子2：这只皮靴号码不小，那只更合适

分词

- 句子1：这只/皮靴/号码/大了，那只/号码/合适。
- 句子2：这只/皮靴/号码/不/小，那只/更/合适。

列出所有词

这只，皮靴，号码，大了，那只，合适，不，小，更

计算词频

- 句子1：这只1，皮靴1，号码2，大了1，那只1，合适1，不0，小0，更0
- 句子2：这只1，皮靴1，号码1，大了0，那只1，合适1，不1，小1，更1

$$\begin{aligned}\cos(\theta) &= \frac{1 \times 1 + 1 \times 1 + 2 \times 1 + 1 \times 0 + 1 \times 1 + 1 \times 1 + 0 \times 1 + 0 \times 1 + 0 \times 1}{\sqrt{1^2 + 1^2 + 2^2 + 1^2 + 1^2 + 1^2 + 0^2 + 0^2 + 0^2} \times \sqrt{1^2 + 1^2 + 1^2 + 0^2 + 1^2 + 1^2 + 1^2 + 1^2}} \\ &= \frac{6}{\sqrt{7} \times \sqrt{8}} \\ &= 0.81\end{aligned}$$

套公式计算

- 句子A: (1, 1, 2, 1, 1, 1, 0, 0, 0)
- 句子B: (1, 1, 1, 0, 1, 1, 1, 1, 1)

词频向量化

处理流程

- 得到了文本相似度计算的处理流程是：
 - 找出两篇文章的关键词；
 - 每篇文章各取出若干个关键词，合并成一个集合，计算每篇文章对于这个集合中的词的词频
 - 生成两篇文章各自的词频向量；
 - 计算两个向量的余弦相似度，值越大就表示越相似。

Outline

余弦相似度、向量空间模型

TFIDF

【实践】TFIDF

LCS

【实践】LCS

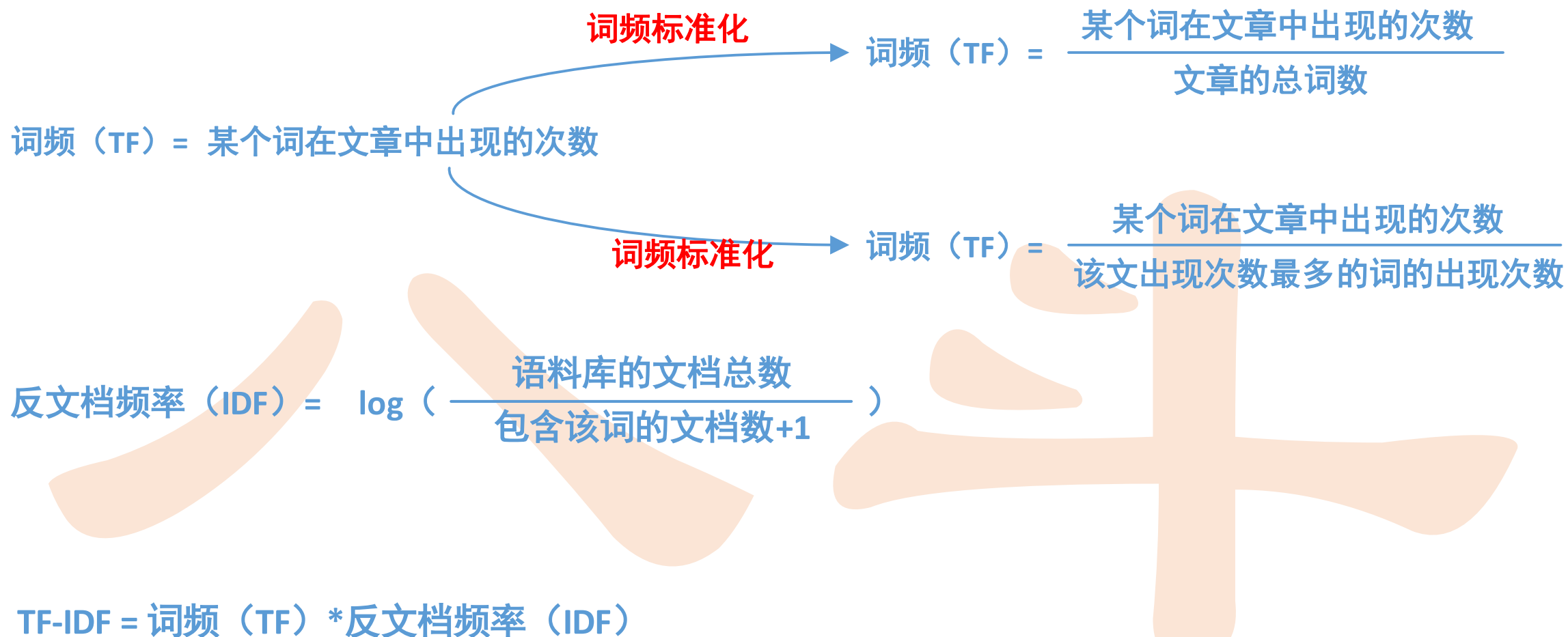
词频——TF

- 假设：如果一个词很重要，应该会在文章中多次出现
- 词频——TF (Term Frequency)：一个词在文章中出现的次数
- **也不是绝对的！** 出现次数最多的是“的” “是” “在”，这类最常用的词，叫做**停用词 (stop words)**
- 停用词对结果毫无帮助，必须过滤掉的词
- 过滤掉停用词后就一定能接近问题么？
- 进一步调整假设：如果某个词比较少见，但是它在这篇文章中多次出现，那么它很可能反映了这篇文章的特性，正是我们所需要的关键词

反文档频率——IDF

- 在词频的基础上，赋予每一个词的权重，进一步体现该词的重要性，
- 最常见的词（“的”、“是”、“在”）给予最小的权重
- 较常见的词（“国内”、“中国”、“报道”）给予较小的权重
- 较少见的词（“养殖”、“维基”、“涨停”）
- 将TF和IDF进行相乘，就得到了一个词的TF-IDF值，某个词对文章重要性越高，该值越大，于是排在前面的几个词，就是这篇文章的关键词。

计算步骤



TF-IDF与一个词在文档中的出现次数成正比，与包含该词的文档数成反比。

相似文章

- 使用TF-IDF算法，找出两篇文章的关键词；
- 每篇文章各取出若干个关键词（比如20个），合并成一个集合，计算每篇文章对于这个集合中的词的词频（为了避免文章长度的差异，可以使用相对词频）；
- 生成两篇文章各自的词频向量；
- 计算两个向量的余弦相似度，值越大就表示越相似。

自动摘要

- 文章的信息都包含在句子中，有些句子包含的信息多，有些句子包含的信息少。“自动摘要”就是要找出那些包含信息最多的句子。
- 句子的信息量用“关键词”来衡量。如果包含的关键词越多，就说明这个句子越重要。
- 只要关键词之间的距离小于“门槛值”，它们就被认为处于同一个簇之中，如果两个关键词之间有5个以上的其他词，就可以把这两个关键词分在两个簇。

- 下一步，对于每个簇，都计算它的重要性分值。

$$\text{簇的重要性} = \frac{(\text{包含的关键词数量})^2}{\text{簇的长度}}$$

- **例如：**其中的簇一共有7个词，其中4个是关键词。因此，它的重要性分值等于 $(4 \times 4) / 7 = 2.3$

- 简化：不再区分“簇”，只考虑句子包含的关键词。下面就是一个例子（采用伪码表示），只考虑关键词首先出现的句子。

八斗大数据内部资料，盗版必究——

自动摘要

```
Summarizer(originalText, maxSummarySize):

    // 计算原始文本的词频，生成一个数组，比如[(10,'the'), (3,'language'), (8,'code')...]
    wordFrequencies = getWordCounts(originalText)

    // 过滤掉停用词，数组变成[(3, 'language'), (8, 'code')...]
    contentWordFrequencies = filterStopWords(wordFrequencies)

    // 按照词频进行排序，数组变成['code', 'language'...]
    contentWordsSortbyFreq = sortByFreqThenDropFreq(contentWordFrequencies)

    // 将文章分成句子
    sentences = getSentence(originalText)

    // 选择关键词首先出现的句子
    setSummarySentences = {}
    foreach word in contentWordsSortbyFreq:
        firstMatchingSentence = search(sentences, word)
        setSummarySentences.add(firstMatchingSentence)
        if setSummarySentences.size() = maxSummarySize:
            break

    // 将选中的句子按照出现顺序，组成摘要
    summary = ""
    foreach sentence in sentences:
        if sentence in setSummarySentences:
            summary = summary + " " + sentence

    return summary
```

总结

- 优点：简单快速，结果比较符合实际情况
- 缺点：单纯以“词频”做衡量标准，不够全面，有时重要的词可能出现的次数并不多
 - 这种算法无法体现词的位置信息，出现位置靠前的词与出现位置靠后的词，都被视为重要性相同，这是不正确的。（一种解决方法是，对全文的第一段和每一段的第一句话，给予较大的权重。）

Outline

余弦相似度、向量空间模型

TFIDF

【实践】TFIDF

LCS

【实践】LCS

Outline

余弦相似度、向量空间模型

TFIDF

【实践】TFIDF

LCS

【实践】LCS

LCS 定义

- 最长公共子序列 (Longest Common Subsequence)
- 一个序列S任意删除若干个字符得到的新序列T, 则T叫做S的子序列
- 两个序列X和Y的公共子序列中, 长度最长的那个, 定义为X和Y的最长公共子序列
 - 字符串12455与245576的最长公共子序列为2455
 - 字符串acdfg与adfc的最长公共子序列为adf
- 注意区别最长公共子串 (Longest Common Substring)
 - 最长公共子串要求连接

LCS作用

- 求两个序列中最长的公共子序列算法
 - 生物学家常利用该算法进行基因序列比对，以推测序列的结构、功能和演化过程。
- 描述两段文字之间的“相似度”
 - 辨别抄袭，对一段文字进行修改之后，计算改动前后文字的最长公共子序列，将除此子序列外的部分提取出来，该方法判断修改的部分

求解——暴力穷举法

- 假定字符串X, Y的长度分别为m, n;
- X的一个子序列即下标序列{1,2,, m}严格递增子序列, 因此, X共有 2^m 个不同子序列; 同理, Y有 2^n 个不同子序列;
- 穷举搜索法时间复杂度 $O(2^m * 2^n)$;
- 对X的每一个子序列, 检查它是否也是Y的子序列, 从而确定它是否为X和Y的公共子序列, 并且在检查过程中选出最长的公共子序列;
- 复杂度高, 不可用!

求解——动态规划法

- 字符串X, 长度为m, 从1开始数;
- 字符串Y, 长度为n, 从1开始数;
- $X_i = \langle x_1, \dots, x_i \rangle$ 即X序列的前i个字符 ($1 \leq i \leq m$) (X_i 计作 “字符串X的i前缀”)
- $Y_j = \langle y_1, \dots, y_j \rangle$ 即Y序列的前j个字符 ($1 \leq j \leq n$) (Y_j 计作 “字符串Y的j前缀”)
- $LCS(X, Y)$ 为字符串X和Y的最长公共子序列, 即为 $Z = \langle z_1, \dots, z_k \rangle$
- 如果 $x_m = y_n$ (最后一个字符相同), 则: X_m 与 Y_n 的最长公共子序列 Z_k 的最后一个字符必定为 $x_m (= y_n)$
- $Z_k = x_m = y_n$
- $LCS(X_m, Y_n) = LCS(X_{m-1}, Y_{n-1}) + x_m$

求解——动态规划法

如果 $x_m = y_m$

	1	2	3	4	5	6	7
X	B	D	C	A	B	A	
Y	A	B	C	B	D	A	B

- 对于上面的字符串X和Y:
- $x_3 = y_3 = 'C'$ 则有: $\text{LCS}(\text{BDC}, \text{ABC}) = \text{LCS}(\text{BD}, \text{AB}) + 'C'$
- $x_5 = y_4 = 'B'$ 则有: $\text{LCS}(\text{BDCA}, \text{ABC}) = \text{LCS}(\text{BDCA}, \text{ABC}) + 'B'$

求解——动态规划法

- 如果 $x_m \neq y_n$, 则 $LCS(X_m, Y_n) = LCS(X_{m-1}, Y_n)$, 或者 $LCS(X_m, Y_n) = LCS(X_m, Y_{n-1})$
- 即 $LCS(X_m, Y_n) = \max\{LCS(X_{m-1}, Y_n), LCS(X_m, Y_{n-1})\}$

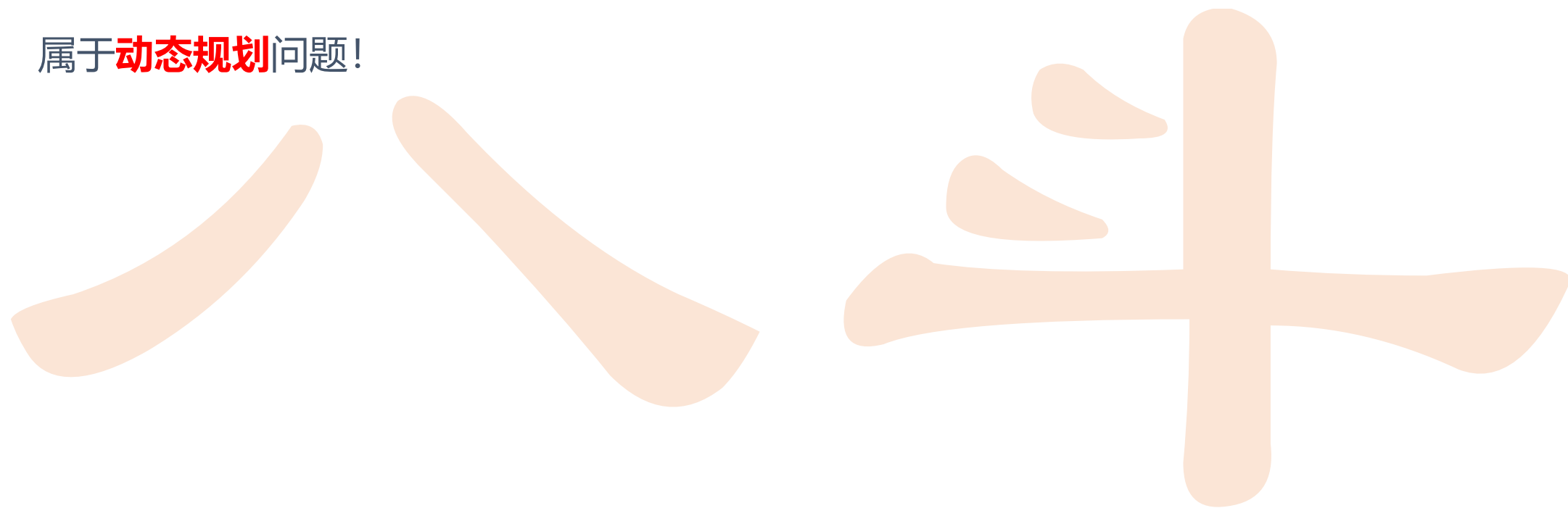
	1	2	3	4	5	6	7
X	B	D	C	A	B	A	
Y	A	B	C	B	D	A	B

- 对于上面的字符串X和Y:
- $x_2 \neq y_2$ 则有: $LCS(BD, AB) = \max\{LCS(BD, A), LCS(B, AB)\}$
- $x_4 \neq y_5$ 则有: $LCS(BDCA, ABCBD) = \max\{LCS(BDCA, ABCB), LCS(BDC, ABCBD)\}$

LCS 总结分析

$$LCS(X_m, Y_n) = \begin{cases} LCS(X_{m-1}, Y_{n-1}) + x_m & \text{当 } x_m = y_n \\ \max\{LCS(X_{m-1}, Y_n), LCS(X_m, Y_{n-1})\} & \text{当 } x_m \neq y_n \end{cases}$$

- 属于**动态规划**问题!



数据结构——二维数组

- 使用二维数组C[m,n]
- C[i,j]记录序列 x_i 和 y_j 的最长公共子序列的**长度**
 - 当 $i=0$ 或 $j=0$ 时, **空序列**是 x_i 和 y_j 的最长公共子序列, 故 $C[i,j]=0$

$$c(i, j) = \begin{cases} 0 & \text{当 } i = 0 \text{ 或者 } j = 0 \\ c(i-1, j-1) + 1 & \text{当 } i > 0, j > 0, \text{ 且 } x_i = y_j \\ \max\{c(i-1, j), c(i, j-1)\} & \text{当 } i > 0, j > 0, \text{ 且 } x_i \neq y_j \end{cases}$$

例子

- $X = \langle A, B, C, B, D, A, B \rangle$
- $Y = \langle B, D, C, A, B, A \rangle$

		j	0	1	2	3	4	5	6
			y_j	B	D	C	A	B	A
i	x_i								
0	x_i		0	0	0	0	0	0	0
1	A		0	↑	↑	↑	↖1	←1	↖1
2	B		0	↖1	←1	←1	↑1	↖2	←2
3	C		0	↑1	↑1	↖2	←2	↑2	↑2
4	B		0	↖1	↑1	↑2	↑2	↖3	←3
5	D		0	↑1	↖2	↑2	↑2	↑3	↑3
6	A		0	↑1	↑2	↑2	↖3	↑3	↖4
7	B		0	↖1	↑2	↑2	↑3	↖4	↑4

Outline

余弦相似度、向量空间模型

TFIDF

【实践】TFIDF

LCS

【实践】LCS

Q & A

@八斗学院
