

与 DF，从两个角度为文本-文档相关性提供了量化依据。例如，当将两篇文章合起来，实现更完整的描述时，一起来考虑  $W_{t,d}$ ， $W_{t,d} = (1 + \log f_{t,d}, d)$ ， $\log \frac{N}{f_{t,d}}$ ，同时，解答了一个问题：什么样的合适度量与相关性有关？

文本性？在文档相关性输出文档的向量，指向陈词滥调的文本和大量，陈词滥调程度越大和大量(5)；**向量模型**，每个文档相关性输出文档的向量，查询时通过文档向量相关性来进行匹配。

用开头词缀的形式来表述加以概括：D：文档表。每个文档提供一个向量，其中每一维对应文档的  $t$ -tf-idf。  
Q：查询表。可操作一个向量，其中每一维对应文档的  $t$ -tf-idf。非完全匹配度  $\beta$ ，R：使用两个向量的相似度来度量查询与文档之间的相关性。向量空间模型的总结。首先，将查询与文档合并成相同的  $t$ -tf-向量空间(也采用其他方法)。其次，计算两个向量的点积(余弦相似度)。最后，按相似度大小排序。将 Top-K 的文档返回给用户。优点：简单直观，可以支持多种不同向量或权重方式。实际效果良好。

缺点：缺乏足够的理解和匹配，同时缺乏  $t$ -idf 值也能造成误导。用  $t$ -idf 来描述文档之间的关系，词项的独立性假设就建立于不成立。例如：贝利、足球 并不等于 足球和贝利 **查询与文档的匹配**，如何计算查询与文档的相关性？最基本的方案：欧氏距离。然而，欧氏距离并不是一个好的选择，它对于向量长度(文档长度)非常敏感。更为合理的方案：余弦相似度。按照向量与查询向量的夹角大小来衡量 **基于余弦的向量相似度更新**，我们使用，用户的查询向量可能无法足够短，需要调整相关性反映现实更多变化。在本质上，这一过程是构造的图表逐步逼近用户输入文档的过程中。如何接近目标文档？可以以目标文档的质心为出发点。某种意义下，由质心的方案，可以以操作质心中的一个点，由于“质心”就是系列点(文档)的质心。我们使用如下方式来计算一簇文档的质心  $\vec{\mu}(C) = \frac{1}{|C|} \sum_{d \in C} \vec{d}$ ，其中，C 是文档的集合，Rocchio 算法提供了一种根据文档信息来归入向量空间模型的办法，其思想在于试图找到一个完美的查询向量，使得足以目标  $q_{opt} = \arg \max_q [\cos(\vec{q}, \vec{\mu}(C_r)) - \cos(\vec{q}, \vec{\mu}(C_{nr}))]$ ，即：使得查询向量尽可能离与相关文档更近，离与之不相关的文档更远。其中，理想情况是，在完美匹配的情况下，不相关文档集合的质心，可以归入以公共权重提供一个完美的查询：  $\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{d \in C_r} \vec{d} - \frac{1}{|C_{nr}|} \sum_{d \in C_{nr}} \vec{d}$ 。

问题在于，我们事实上并不能获得完整的文档/不相关文档集合。Rocchio 算法(1971)实际使用的方案如下： $\vec{q}_{opt} = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{d \in D_r} \vec{d} - \gamma \frac{1}{|D_{nr}|} \sum_{d \in D_{nr}} \vec{d}$ ，其中， $\vec{q}_0$  为初始查询向量。

D<sub>r</sub> 为已知相关文档的向量集合，D<sub>nr</sub> 为已知不相关文档的向量集合。α 为初始查询向量。为根据用户输入文档经验进行校正。由于，新的向量指向符合需求的相关文档向量，指向不相关文档向量 **正反馈** 与 **负反馈**。反映反馈的权重往往大于负反馈。用更更改与查询向量的向量权重，而不是错误权重。相关的，可以通过设置  $\beta > \gamma$  来给予正反馈更大的权重。很多系统甚至只允许正反馈，即  $\gamma = 0$ 。收集相关的正反馈往往比消除 TF-IDF 文档的局限性。使用  $t$ -idf 来求词项频率和逆文档频率且容易理解。但是， $t$ -idf 仅以“词频”量度的重要性。无法将文档提供的信息值以词频与上文的关联：**Edg**，武校时度(“词频”武校)：因此，我们需要一种能够表示出上文与下文的词频表示方法以**描述词频与上文的关联**。如何构造词频量度以及目标。获得词频与上文的关联关系表，Word2vec 模型实现：两两计算矩阵。根据上文预测中(CBOW)，根据中词预测上文(Skip-gram) **Word2vec 模型的基本流程**。我们以 CBOW 模型为例，介绍其工作原理(Skip-gram 模型基本思想类似)。CBOW 本质上是在获得上述的 O(1) 中词的过程中，介入了词频这种映射，由于：深度学习中的 CBOW 算法详细步骤如下：使用 one-hot 的 O(1) 向量表示词频。计算结果值(深度学习基础)。计算预测值的中间值的概率： $P(\text{某词} | \text{某词频}, \text{某词频}, \text{某词频})$ 。计算结果值(深度学习基础)。有了预测值的上文关系表。无监督。通常来说，可适用于各种 NLP 任务。缺点：无法解决一词多义的问题。例如 play 单词和 play football，Word2vec 是一种静态的表示，其词向量表示一旦训练确定就不会再更改。而且，虽然通用性强，但是无法针对特定任务进行动态变化**三常见的排序算法**。Pointwise：对相似性为类或等级的问题；输出：两两对立的分类(有/无)，训练或没有有训练；Pairwise：比较一对网页的相似度，e.g. 相关 / 不相关；输出：两两对立的回归关系；Listwise：对整个数据集进行建模，输出：整个集合的完整排序，往往使用带有特定排序**Pointwise 类排序算法**。给定数据：训练集中的任意一个文档-词对，都可映射到一个值或一个有界的类别(如良/差)。组成的，第一个文档-词对，Pointwise TR 将训练数据分类/鉴别)。常见的模型类型包括：回归，将每个文档-词对映射到具体数。分类，将问题转化为一方面具有两类的结果(多/少/多/少)。有序回归，在数据映射到值的的同时保持排序的值的有序性。Pointwise 类排序算法可以单独使用，也可以广泛地应用于回归、分类算法，然而，其局限性也较为明显。首先，Pointwise 类方法往往难以成为文档相关性问题的首选，而不在于文档相关性问题的天然特性，Pairwise 类方

为了解决这一问题提供了新的方案。其次，不同页面所对应的文档，无论其篇幅与质量如何不同，对关联数据的出现，都基于同一程度上影响了新数据 **Paisirwe 类排序算法**；基本思路：将网页问题转化为类似为问题（二）或三的问题，每比较一个查询访问两次文档，衡量两个文档的偏序(Partial Order)，分类的目的不在于原文档链接是否有效，讨论的是对于  $i = 1, \dots, N$  ( $N$  表示所有文档数)而言，相对于 Pointwise Pairwise Ranking, Paisirwe 类排序法通过调整排名本来的顺序关系，实现了从绝对相关性(分值)到相对相关性的转化。然而，Paisirwe 类排序法具有自己的缺陷：首先，两两对比导致计算成本大为提升，并需要开发复杂的 Paisirwe 类算法；虽然它自然存在不平衡问题的影响，最后，Paisirwe 类算法无法完全呈现全局信息的合理性，由此导致了后者的缺点：**Listwise 类排序法**，局限性：样本不平衡的影响，如相同 Query 的 Doc 数量相差很大，从最坏的情况说，Doc (即) 数量较多的 Query 将其相关的 Doc，从而产生许多 Doc 可通过交叉验证加以缓解，而例性：算法无文本信息局部性，某些文档结构排错，发生在非多位元 Paisirwe 条件中是可行的；局限性：7 篇文档，相关度从 1 到 3，理想情况下的排列为 3 2 2 1 1 1 1，如有如下两种排列方式：Rank A: 3 2 2 1 1 1 Rank B: 3 2 2 1 1 1，这两个排列，在 Paisirwe 条件下是等价的，然而，前者显然更符合，如前述，用住往关心核心在前靠位值的文档 Listwise 类排序法；基本思想：遵循向整体转移信息进行优化，因此直接将排列完整化作为学习的目标，通常情况下，这一项需要用至少两种方式，直接采用某种 IR 指标进行评价和选择，直接按上述方法得到的损失函数通常情形下，由于全球化的应用，Listwise 类排序法可以取得相当于 Pointwise 和 Paisirwe 类排序好的效果，然而，Listwise 类用法也会面临一些小的挑战，例如两个网页并列的情况，相当于 Paisirwe 中文档的相关度的否定的情况 **Paisirwe 类的改进型**，在网页操作之间，网页间的数据提供一条边，形成大量的有向图，PageRank 算法的基本思想：优质网页的 PageRank 值高，→ 优质网页指向低质量的网页，输入链路强度，在链路上的受权程度越高 输入链路的强弱是否指推荐度，页本身是否具有权威性，则根据页面的链接数：被选中的一级，体现源网页是否显著推荐 **PageRank 的计算公式**， $\text{PageRank}$  的公式定义为：
$$PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} PR(p_j) \frac{PR(p_j)}{\sum L(p_j)}$$
$$\text{为: } a = dR + (\frac{1-N}{d}) ee^T, ee^T \text{ 是全 1 矩阵、} PR(pi) \text{ 为网页 pi 的 PageRank、}$$
$$(PR(pi)) \text{ 为指向网页 p_i 的所有网页 p_j 的 PageRank 乘积、} L(p_j) \text{ 为网页 p_j 发出的链接数量、}$$
$$\text{为阻尼系数，取值在 } 0-1 \text{ 之间、} N \text{ 为网页总数、} M(p_i) \text{ 为输入 p}_i \text{ 的网页集合、} R \text{ 为由统一初始权重、}$$
$$\text{为邻居矩阵转置、} P \text{ 为网页之间的邻接矩阵、} \text{PageRank} \text{ 的迭代过程是：} \text{求得初始迭代计算结果、}$$
$$\text{得到 PageRank、} \text{重复给每个网页赋值、} \text{例如 } i=1, \dots, n \text{ 个网页，依次求取 PageRank、}$$
$$\text{收敛标准：前 } n \text{ 次迭代中的两次连续 PageRank 的差值小于阈值，则认为收敛、}$$
$$\text{近似 PageRank 算法中的两个关键因素：} \text{PageRank 的计算方式，使用了一个马尔可夫链模型，如果出、}$$
$$\text{现故障节点点线性分布，如何处理？、网络特点：只有一条自回路向外，没有他处出口，截止节点，没、}$$
$$\text{何发出，如同漏网，此外，孤立节点的存也在客观上会产生一定影响，及没有任何入口，单独使用马尔可夫链无法统计、}$$
$$\text{有初期概率，不再更新，也不影响其他节点几点特殊情况的处理：Restart 机制 - 回顾 PageRank 的历史、}$$
$$\text{其中的 } (1-d)/N \text{ 部分，相当于一给定概率重新选起点，此时，所有节点在一个给定概率率中成为公式、}$$
$$\text{由此，就消除了陷阱和网络上的干扰 (所有节点平等连接)、} d \rightarrow \text{一般选取为 } 0.85 \text{ 左右 (Google 的设计师把 PageRa、}$$
$$\text{k 中的收敛性说明，如前述，PageRank 的计数过程，实际上是一个马尔可夫过程，取 } A = dM + (1-d)/Nee^T, \text{其中 } M \text{ 为转移矩阵，为所有元素大于 0 的列向量，则有 } Pn+1 = AP_n, \text{形成的值为收敛收、}$$
$$\text{敛，论述之中从理论上证明，不讨论 PageRank 值何时收敛，这种方法能解释了网页排名中的正反馈效应、}$$
$$\text{的真实性，马尔可夫方程的三个必要条件：} \text{转移矩阵 } A \text{ 为非负实矩阵、} A \text{ 矩阵所有元素都大于 } 0, \text{并且每一列之和等于 } 1, \text{转移矩阵 } A \text{ 不为零的，当满足强连通时，} A \text{ 可为不可逆，而 Restart 价、}$$
$$\text{了这一条件，转移矩阵 } A \text{ 为非奇异的可，这三个条件，PageRank 算法的收敛性，因此保障了其收敛性、}$$
$$\text{HITS 的两个基本原则：} \text{基本原则是权威性和相关性、基本原则 1 好的 Authority 会很多很多的 Hub、}$$
$$\text{问，权威级 2：好的 Hub 会指向多个好的 Authority，因此，在 HITS 算理中，每隔一段时间要、}$$
$$\text{两个权威做 HITS 更新的计算过程，基本的先前的假设，HITS 的对应原理上：HITS 归一化迭代方、}$$
$$a_k k+1 = \frac{MT}{|M|h_k}, h_k k+1 = \frac{Ma_k k+1}{|Ma_k k+1|}$$
$$\text{取初值 Hub, Authority 均为全 1 向量、}$$
$$\text{首先，根据关系矩阵与初始值最大的少数页面，以及这些页面有链接关系的页面，作为待选集，此时，对所有、}$$
$$\text{点的 } a(k) \text{ 与 } h(p) \text{ 进行标准化处理，可都变为 1，然后，迭代执行两个步骤，即基本算理所对应的两个公式，才、}$$
$$\text{这一步，直到达到收敛为止，注意每一步中的一致化处理！输出 Authority 或 Hub 用户的页面(即) } \alpha \text{ 个性化推、}$$
$$\text{荐是基于该用户的偏好：} \text{基本想法：用的一般相似性稳定，因而，用户偏好他们的偏好(她、}$$

$\text{WF}_t$ , 从两个维度与向量文档相关性进行了量化优化。例如, 当将两篇文档合并, 实现更紧密的链接时, 乘起来得到  $D_t = (1 + \log t f_{t,d}) \cdot \frac{\log N}{f_{t,d}}$ 。同时, 解答了一个问题: 什么样的合适度量衡量了相关性?

关联? 在少数文档内多次出现词语, 词频越高则相关性越大而大量文档(5); 文档空间稠密? 每个文档包含词语出现的次数越多, 查询时遇到该文档向量的相关性越低; 或者, 当开头单词形成的初始表示添加以权重 $f_{t,d}$ , 文档向量越短, 每个文档指向一个向量, 其中每一维对应词语的  $t\text{-idf}$ 。

Q: 查量级表, 可视为一个矩阵, 其中每一维对应词语的  $t\text{-idf}$  值。 $P_r$  完全匹配集合  $R_r$ ;  $R_r$  使用所有文档之间的相似度来衡量文档与查询之间的相关性, 向量空间的距离的总结。接着, 将文档与查询向量对应的  $t\text{-idf}$  取最大值(也采用其他方法)。其次, 计算两个文档集的某种相似性(如按欧几里得), 最后, 按相似度大小排序, 将 Top- $K$  的文档返回给用户。 $t\text{-idf}$ : 简单地说, 可以支持多种不同度量或权重方式, 用效果最差的方式来描述; 缺乏足够的理解和知识时, 同时使用  $t\text{-idf}$  值也可造成误导。用户无法通过文档之间的关系、词语的独立地位或独立性上建立。例如: 贝利、足球和排球 + 排击是查询的文档。如何计算文档与查询之间的相似度? 最基本的方案: 欧氏距离。然而, 欧氏距离并不是一个好的选择, 它对于向量长度(文档长度)非常敏感; 更为合理的方案: 余弦相似性。按照向量与查询向量的夹角大小来度量基于汉字的向量更新。下面我们提到, 用户对查询文档可能无法一致地回答, 需要调整相关度反映现实世界中数据。在本质上, 这一过程是查询图形的逐步逼近并逼近用户文档结构的过程。例如近似接近标签文? 可以从目标文档的中心点出发。某种情况下, 由向量表的中心点, 可以视作待处理过程中的一个点。由于, “质心”就是—系列点(文档)的均值。我们使用如下方式来计算一类文档的质心  $\vec{\mu}(C) = \frac{1}{|C|} \sum d \in C$  其中,  $C$  是文档的集合,  $|Roch|$

公式提供了这样一种相关度反馈信息即回到向量空间中搜索的点。其思想在于试图找到一个完美的最佳使用文档, 满足以下目标  $\hat{q}_{opt} = \arg \max [\cos(\vec{q}, \vec{\mu}(C_r)) - \cos(\vec{q}, \vec{\mu}(C_{nr}))]$  找到: 使得差值尽可能高与查询文档一致, 高于之不相近的文档更优。其中, 理想情况是, 在可知完整的相关/不相关文档表格的情况下, 可以得到  $\hat{q}_{opt} = \frac{1}{|C_r|} \sum \vec{d}_j \in C_r \frac{\vec{d}_j}{|C_{nr}|} \sum \vec{d}_j \in C_{nr}$

问题在于, 我们事实上并不能获得完整的文档/不相关文档集合。Rocchio 算法(1971)实际采用的方法如下:  $\tilde{q}_m = \alpha \tilde{q}_0 + \beta \frac{1}{|D_r|} \sum \vec{d}_j \in D_r \frac{\vec{d}_j}{|\gamma| \frac{1}{|D_{nr}|}} \sum \vec{d}_j \in D_{nr} \frac{\vec{d}_j}{|\gamma|}$

中,  $D_r$  为已知相关文档的向量集,  $D_{nr}$  为未知不相关文档的向量集合。Q 为初始查询向量,  $\alpha, \beta$  为权衡参数,  $\gamma$  调节收敛速度。由此, 新的向量由不相关文档到相关文档的移动, 远离不相关文档的最近正交投影位置。反方向的移动往往比大负向量, 用系统更大符合需求的推荐策略, 而不是错误策略。相似的, 可以通过调整  $\alpha, \beta$  来给予不同重要度的文档。很多系统甚至只允许正反馈, 即  $\alpha = 0$ 。收集更多的正反馈往往优于  $-IDF$  特征的非线性性。使用  $t\text{-idf}$  来表示词语频率快速且容易理解。但是,  $t\text{-idf}$  仅以“词频”作为标准, 无法体现词语的重要性差异以及词语的上下文。因此, “语义场”、“语境分析”等。因此, 本文首先展示一下上文定义的  $t\text{-idf}$  表示的方法及如何改进提高对上下文的关系。改进后的设计以提高为目标。获得词语间上下义关系的关键要素: Word2vec 模型训练上的改进, 再结合上述原理。根据上文案例中的(CBOW), 根据中心词预测上下文(Skip-gram) Word2vec 模型的基本流程。我们基于 CBOW 模型, 介绍其具体工作原理(Skip-gram 模型基本思路略过), CBOW 本质上获得上下文向量 → 中心词的表示, 如何获取特征的映射? 基于深度学习工作的 CBOW 详细步骤如下: 使用 one-hot (0/1) 向量表示的词, 计算编码层(深度卷积核); 计算谓词的中间词的概率; 计算交叉熵损失, 反向传播优化模型; 多个不同词源窗口进行训练测试; 最终得到各个谓词的权重: Wx1 Word2vec 的优点: 快速, 有效; 解决了词项之间的上下文关系。无监督, 通用性强, 可适用于各种 NLP 任务。缺点: 无法捕捉太多次数的词语, 例如 play music and play football, Word2vec 是一种神经网络方式, 词项被逐一训练但训练就不会更快一些, 虽然通常很快, 但是无法针对特定任务动态调整这是三次常见的排序学习算法。Pointwise 将退化为了分类或回归问题。输出: 网页对不同的行为(有/无), 归因或没有序回归值。Pairwise: 比较一对网页的相关性(是否、相关、不相关。输出: 网页对不同的偏好关系。Listwise: 对整个页面集进行比较并输出: 整个集合的完整排序, 往往在依特定顺序排序 Pointwise 类排序算法。基本原则: 训练集中的任何一列  $H_i$  行, 都可映射到一个正值或一个负值的类型(如良/恶)。组成的, 给定一个查询文档, Pointwise LR 将视图特征(类别/级别)。常见的数据类型包括: 回归, 回归, 将每类文档输入到线性分数。分类, 将问题转化为一面值为有序列表(多类/多分类)。有返回值, 在数据集到模型的训练保持样本中的有关系的元素。Pairwise 类算法具有可以简单且广泛地使用回归/分类法。然而, 其局限性较为明显, 因为 Pairwise 类算法往往更注重于文档间的相互影响, 而不注重文档与文档间的长期关联性; Pairwise 类方法应用, 为解决这一问题提供了新的方案。其次, 不可忽略查询文档的长尾, 尤其相关文档数量不同, 对误差敏感的程度不相同。一定程度上影响 Pairwise 类算法。基本原则: 如何将排序问题转化为二元问题(二分类问题)? 类似比较每一个影响因素 Pairwise 文档, 衡量两两文档的相似性(Partial order), 分目的在于将三元文档相关性转换为二元的, 即对每对文档  $(p, q)$ ,  $0 \leq p, q \leq 1$  (表示两两文档相似程度)。相对于 Pairwise 类算法, Pairwise 类算法法具有直接构造两两文档间的偏好关系, 实现了从相对偏好(三元)到绝对偏好的自然转换。然而, Pairwise 类算法法具备良好的可扩展性。首先, 两两对比对数随样本增大而增加, 并计算复杂度随之增加。Pairwise 类算法法天然存在不平衡问题的影响。最后, Pairwise 类算法法无法完全解释合理性。由此引出了最后一类算法: Listwise 类算法。局限性: 样本不平衡的问题。如不同 Query 的 Doc 数量相差很大, 从偏见的角度来看, Doc (Query) 数量较多的 Query 将掩盖其他的 Query, 从而产生偏差; 计算时可通过加权等方式加以缓解。缺点: 无法体现非全局性偏好, 某个成文文档排序错误, 发生在不平衡位置。Pairwise 条件件是平等的。例如: 7 篇文章, 相关从 1 到 3, 理论最优的排序为 3 2 2 1 1 1 1。如有知两个样本, RankA: 2 3 2 1 1 1 1 Rank B: 3 2 1 2 1 1 1。这两个排序, 在 Pairwise 类算法上是平等的, 然而, 后者更符合现实意义。如前所述, 用户往往关注其在喜欢位置的文档 Listwise 类算法。基本原则: 直接将向量特征结果进行迭代优化。直接将其推荐给团队用于学习的实验。通常做法下, 这一团队采用以上两种策略。直接使用某种 LR 指标将排序结果与期望。直接验证面向页面的推荐效果通常情况下, 由于不良优化的作用, 使用 Listwise 类算法排序可以取得相较于 Pointwise 和 Pairwise 类算法更好的效果。然而, Listwise 类算法也有面临的一些小的缺陷, 例如单个网页并列的情况。相当于 Pairwise 中文档中标签为 0 的情况 PageRank 的设计理念。将网页作为一个点, 网页间的超链接构成一条边, 形成巨大的有向图。PageRank 背后的基本思想: 网页质量是由其所接收的、一定还来自高质量、三重激励: 链入链接数; 单纯靠文字上的受欢迎程度; 网页用户本身是否推荐者。因此要评估是否具有权威性, 链入链接页面的链接数: 选中中的几率。资源网页是否是受推荐 PageRank 的计算方法。PageRank 的公式如下:  $PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j)}{L(p_j)}$  PageRank 转移概率矩阵为:  $A = dR + (\frac{1-d}{N}) ee^T$ ,  $ee^T$  是一矩阵。PR(pi) 为网页 pi 的 PageRank

• PR(pi) 为指向网页 pi 的所有网页 p<sub>j</sub> 的 PageRank。L(pi) 为网页 pi 发出的链接数量; 为阻尼系数, 取值在 0-1 之间。M(pi) 为输入 pi 的网页集合。R 为同一化链接矩阵。M 为网络链接矩阵 PageRank 的实现过程。PageRank 的迭代计算方法。采用迭代化的链接矩阵 PageRank 值。首先给每个网页赋予一个初值, 例如 1/N。然后, 利用上一次的值进行迭代更新每次计算, 近似收敛 PageRank 中的两种特殊情况。例如 Rank 1 计算后, 实际上是一个马尔科夫链过程。如果出现循环或终止性节点, 如何处理? 陷落节点: 只有一条指向自己的边, 没有 outgoing 边。终止节点: 没有 outgoing, 如同漏网。此外, 孤立节点的存在也会导致产生一些问题。没有任何入边, 单独使用马尔科夫链无法提供有效的初始化概率, 不再更新, 也不影响其他节点(及其特指性的解决: Restart 机制。回顾 PageRank 的公式中的 (1-d)/N 的节点, 相当于一以概率率重新选取节点。因此, 每节节点以一固定概率率选中作为新节点。由此, 防止了陷落和漏网的干扰(所有节点全连接), d — 取平均值 0.85 左右(Google 的设计选择) PageRank 中的收敛性问题。如前述, PageRank 的计算过程, 实际上是求一个马尔科夫过程。取 A = dM + (1-d)/N e e<sup>T</sup>, 其中 M 为转移矩阵。为所有元素都大于 0 的列向量, 则有 Pn+1 = APn, 形成马尔科夫链。[Note] 已从图上提出来, 不过, 仍讨论 PageRank 相关问题吧, 这种算法保证了网页间的良性收敛性。马科夫链的三个必要条件: 为转移的。A 为马尔科夫转移矩阵, A 矩阵所有元素都大于 0, 其任意一行之和等于 1, 称为双随机。为非常规的。这条性质, PageRank 算法满足, 因此确保了马尔科夫链 HITTS 的两个基本原则: 基本假设与核心概念相互呼应。基本假设 1: 好的 Authority 会很多很多的 Hub 网络。基本假设 2: 好的 Hubs 会指向很多个好的 Authority。因此, 在 HITTS 算术中, 每个网页需要两个权威值 HITTS 算法的计算过程。基于先前的基础假设, HITTS 的计算过程如下: HITTS 目前一代迭代方法  $a_k+1 = \frac{M^T h_k}{|M^T h_k|}$ ,  $h_k+1 = \frac{M a_k+1}{|M a_k+1|}$  取初始值  $h_0$ , Authority 为全 1

首先, 根据统计方法与查询最相关的少数页面, 及这些页面有链接关系的页面, 作为特殊集合, 此时, 对所有页的  $a(k)$  与  $h(p)$  进行初始化, 可都设为 1。然后, 迭代计算两个步骤, 即基本假设所对应的两个公式, 直到收敛。直到收敛为止, 注意每一步中的一致化过程! 输出 Authority 或 Hub 值较高的页面(即个性化推荐。基于特殊的推荐: 搜索引擎。基本想法: 用的一般的相似性稳定。因此, 指出, 指出一种特殊化, 他们



