核心概念: 1) 权威网页 (Authority): 指某个领域或某个话题相关的高质量网页 2) 中心网页 (Hub): 类似中介, 指向 HITS 目的: 在海量网页中找到并区分这些与用户音询主题相关的高质量 "Authority" 5 "Hub" 网页,尤其是 "Authority" (因为更能满足用户的信息需求)

$$\forall p, a(p) = \sum_{i=1}^n h(i)$$

$$\forall p, h(p) = \sum_{i=1}^{n} a(i)$$

假定邻接矩阵为M, Authority向量为a, Hub向量为h

则有如下迭代式:

 $a_{h+1} = M^T M a_h h_{h+1} = M M^T h_h$ 解特征方程 (A-λl) x=0 且Σx<sub>i</sub>2=1 得機  $\bullet \quad a_{k+1}=M^Th_k, \quad h_{k+1}=Ma_{k+1}$ 

或者,可采用如下迭代式:

- $a_{k+1} = (M^T M)^k M^T a_0$ ,  $h_{k+1} = (M M^T)^{k+1} h_0$
- 其中、a<sub>0</sub> h<sub>0</sub>为Authority/Hub向量的初始值、可设为全1向量

优缺点: 相比于 PageRank,HITS 算法是一种能够区分网页功能的排序算法 优点: 1) 更好地描述互联网组合特点 2 主题相关,因此可以单独用于网页排序 缺点: 1)需要在线计算,时间代价较大 2)对链接结构变化敏感,且依然可能受 到"链接作帐"的影响 3) 初始市面的选择对音询结果有影响 需要找到好的选择管法

基于内容的推荐技术 (用户的偏好一般相对稳定): 物品画像&用户画像(基于用户评分进行加权)采用余弦相似性度量进行? 分 优缺点: 优点: 1)每个人的推荐过程相互独立,不需要其他用户的数据 2)可以为具有独特偏好的用户进行有效指 荐 3) 可以推荐新项目或非热门项目 推荐结果有着较好的可解释性 缺点: 1) 找到合适的特征是一件困难的事 2) 给新用 户推荐项目,永远是一个困难的任务 3) 过度特化 (Overspecialization) 现象

多样化评估:最大边界相关性 MMR 基于路径推荐:知识图谱代替向量化画像,基于图谱上的游走实现推荐, 偏见 Bias 问题: 位置(第一更受关注), 模态(与众不同的模态吸引关注), 关键词效应 双向选择问题: 用户<->物品, 稳定匹配

<mark>办同过途</mark>算法实现:1)基于内存(Memory-based)的协同过滤(基于现有数据与简单度量运算进行推荐,可进一步细分 为基于用户(User-based)与基于项目(Item-based)) 2)基于模型(Model-based)的协同过滤(基于现有数据训修 模型,通过模型进行推荐,代表性方法如矩阵分解(Matrix Factorization)、深度学习等)

User-based: (预测时忽略空值): 要寻找最近邻

$$sim(a,b) = \frac{\sum_{p \in product(P)} (r_{a,p} - \overline{r_a}) (r_{b,p} - \overline{r_b})}{\sqrt{\sum_{p \in product(P)} (r_{a,p} - \overline{r_a})^2} \sqrt{\sum_{p \in product(P)} (r_{b,p} - \overline{r_b})^2}}$$
Assesses as for a function of user to the state of t

计算相似度时,如果某个user对某个item未评分。则对应的r直接设为0. 不用减去平均分

$$pred(a, p) = \bar{r_a} + \frac{\sum_{b \in neighbors(n)} sim(a, b) \cdot (r_{b, p} - \bar{r_b})}{\sum_{b \in neighbors(n)} sim(a, b)}$$

$$r_{ix} = rac{\sum_{j \in N(i:x]} \!\! \left[ s_{ij} \right] r_{jx}}{\sum \!\! s_{ij}}$$
 计算方式与User-based的Similarity相似

## 相似度计算时,如果未评分,直接设为 0,不用减平均数。

基于物品的推荐效果更好原因:物品属性单一,用户偏好多样。

基于内存优点:不受多模态非结构化信息表征与特征洗取用扰:

缺点: 冷启动、稀疏性、热度偏差。

冷启动问题: (1)先非个性化推荐 (2)借助个人信息或其它网站浏览信息 (3)诱导式推荐迭代收集用户反馈 (4)与基于内容混 合 (5)Side information: 众包文本、知识图谱

基于模型(Model-based): 潜在因子 •  $\min_{p,0} \sum_{(i,x) \in R} (r_{xi} - q_i \cdot p_x)^2$  这一公式,其本质是一个误差平方和(SSE)。 · 因此,我们的目标正是通过优化这一SSE,以获得潜在因子的估计值

钻阵分解

- 然而,潜在因子的维度 K,也是一个潜在的麻烦 我们需要训练的参数一共(M+N)\*K个,随着K的增长而增长
- 如果K足够的大,在数据稀疏的情况下,可能并没有那么多训练样本
- 因此,过高的 K 可能导致过拟合的问题
- 解决讨高的 K 带来的讨拟合?引入正则项,避免过大参数值

$$\min_{P,Q} \sum_{\text{training}} (r_{xi} - q_i p_x)^2 + \left[ \lambda_1 \sum_{x} \|p_x\|^2 + \lambda_2 \sum_{i} \|q_i\|^2 \right]$$

### 概率矩阵分解: 另一种用路则着银干参数本身的生成规律

- 在概率矩阵分解中、两个参数矩阵均限从均值为0、预设方差的高斯分布 在数据稀疏四有噪声的情况下,是否可能引入某些规律。可以实现对于 参数更好的描述。比如、参数符合高斯分布?  $\rho(U \sigma_2^2) = \prod_{i=1}^d S(U/0, \sigma_2^2 I)$   $\rho(V \sigma_2^2) = \prod_{i=1}^d S(V/0, \sigma_2^2 I)$
- · 开放问题:如何求解PMF中的矩阵参数?基于R与U、V联合分布优化求解 · WOLDERSTEINSTON 假设:: 两个潜在因子矩阵相互独立,各用户 项目业上月分布 在參数完成训练后、即可简单助得到相应的完整评分矩阵
- 概率矩阵分解 (Probabilistic MF)

### $P(R(U, V, \sigma^2)) = \prod_{i=1}^{N} \prod_{i=1}^{N} [N(R_i|U_i^TV_i, \sigma^2)]^{i}$

非负矩阵分解: 对求解过程添加非负约束,以保障元素的非负性 矩阵分解加约束:来加入更多信息与假设,往往可 一种最常见的约束方式: 社交约束 (基于偏好相似的社交约束虽然效果很好, 但未必始终成立)

情境信息:帮助我们理解用户意图的信息。服务提供者可借助情境信息,为用户提供更精确的信息检索和过滤服务。 音询概念: 一组有着相同语义的音询词 (可以解决音询词的稀疏性问题、同时、更规范地解释音询上下文

查询上下文可以帮助更好的理解用户查询词。 上下文感知搜索: 模型准备阶段+感知查询阶段

### 几种常见的查询上下文类型: · 查询重组 · 查询特化 · 查询泛化 · 一般关联 情境感知的推荐任务

推荐系统中的情境:情境是那些在重复讲行同一活动时可能发生变化的因素

获取情境:显示获取 (用户选) 隐式获取 (从应用数据获取) 推断获取 (从用户近期行为推断归纳)

情境筛选:直接询问用户 通过特征选择(主成分分析(PCA)和线性判别分析(LDA)) 通过统计分析(统计测试,如信息增益、5 信息、Freeman-Halton Test 检验等) 情境感知推荐的三种范式:情境预讨波(事先根据情境信息筛选、分隔数据记 录) 情境后过滤(产生推荐结果后,根据情境信息进行筛选) 情境建模 情境信息整合到推荐系统: 传统推荐: Users × Item Ratings 情境感知的推荐: Users × Items × Contexts Ratings

数据预处理: 问题: 1. 数据测量、采集等过程中出现的错误 2. 噪声、离群点、缺失值等数据问题 3. 重复数据 **初外理方式:数据采样 维度旧约 数据索款** 

数据采样: 指选择一部分数据对象的子集进行分析的常用方法 问题: 采样缺乏代表性,将影响对于原数据集的还原程度,进而产

- 影响平柱効果的重要因素之一早平柱的样本容量 最基本的采样技术为简单随机采用 (Simple Random Sampling) 较大的样本容量更能完整代表数据。但降低了采样的收益
- 对于所有对象。采用简单的等概率方式讲行采样 较小的样本容量在采样收益上更高,但可能造成信息的损失
- 无故回采样:被采中的对象从整体中删除,仅可选中一次
- 有放回采样:被选中的对象不会从整体中删除,可再次被选中 通过分组采样的方式,可以近似确定适当的样本容量 两种方式无本质区别, 但有放回采样较为简单(概率不变) 渐进式采样从训练结果的角度确定采样规模 (优点:不需要在
- · 更为复杂的方法包括分层采样 (Stratified Sampling) 开始就确定采样的容量 缺点: 计算开销大 (需要多次迭代)) • 对数据进行分组,从预先指定的组里进行采样

数据归约: 维度归约,可以删除不具有区分度的特征,同时可能降低噪声

方法一: 中成分分析 (PCA): 通过正交变换将一组可能存在相关性的变量线换为一组线性不相关的变量, 装换后的这组变量则中成分

最大特征值对应的特征向量可以最大化投影方差。

完成第一个方向基选择后,第二个方向基应选择使其不存 在相关性,故优化目标为协方差为0,即  $\frac{1}{m}\sum_{i=1}^{m}a_{i}b_{i}=0$ , 等价于协方差矩阵 $C=\frac{1}{m}XX^{T}$ 对 ,即除对角线 (方差) 外的其它元素化为0。故优化 目标转变为寻找矩阵P,满足 $PCP^T$ 是一个对角矩阵。



最大 K 个特征值的特征向量对应的线性组合是 K 个新指标(K 个特征值的比重反映了主成分的信息量,一般>0.85)

输入: n维样本集 $X=(x_1,x_2,\cdots,x_m)$ , 要降维到n' PCA 特点(1)依赖原始变量也只能反映原始变量(2)PCA 内在假设 输出: 降维后的样本集Y 之一是原始变量直接存在一定关键。 若原始变量相互独立则降维生

败,数据越相关,降维效果越好(3)PCA 的结果未必清晰可解释 (1) 对所有样本进行中心化 $x_i=x_i-rac{1}{\mathbf{n}}\sum_{j=1}^{\mathbf{n}}x_j$  差异主要体现在描述内容、结构和表示、应用领域和时间维度上 (2)计算样本的协方差矩阵C=  $XX^T$  **ゅっ**し

(3)求出协方差矩阵的特征值及其对应的特征向量 5 (4)将特征向量按特征值大小从上到下按行排列成矩阵 有内面 有内面 有内面 《实体、属性 医小术实体、关系、实体、三元组 《事件、论元集会、逻辑关系》多元组 取前k行组成矩 以前k行组成矩阵 (5)Y = PX即为降维到k维后的数据

方法二: 特征子集选择 (而不是归纳新特征) 去除冗余特征和不相关特征, 或为特征赋予不同权值 数据离散 (将连续属性变换为分类属性): (1)二元化: 将连续或离散属性转化为多个二元属性 0/1 (2)非监督离散化: 不用类别信息 (等宽/等频率/等深) (3)有监督离散化 (基于熵)

### 知识图谱概述

<mark>言息抽取</mark>含义:从语料中抽取指定的事件、事实等信息, 形成结构化的数据 信息抽取是整合与分析的基础 信息抽取与信息检索: 两者密切相关, 却又存在鲜明差异 1) 功能不同 · 检索: 从文档集合中找文档子集 · 抽取: 从文本中获取用户 感兴趣的事实信息 2) 处理技术不同 · 检索: 可以采用统计与关键词等技术 · 抽取: 需要借助自然语言处理技术 3) 使用领域不同

信息抽取基本任务:命名实体 NE (最重要)、模板元素 TE (模板元素又称为实体的属性,目的在于更加清楚、完整地描述命名实体) 共指关系 CR(如果不同的命名实体表达了相同的含义,即为共指关系,也称为等价概念。共指关系的抽取任务在于抽取关于共指表达 的信息)、模板关系 TR (实体之间的各种关系,又称为事实。通过关系抽取,将实体关联起来,并为推理奠定基础)、背景模板 ST (又

知识图谱表示:由结点和结点之间的边组成。 结点表示概念(或字体)。 边表示关系(或属性)。 • 在数学上,知识图谱表现为一个有 一般而言,知识图谱中的节点用来表示概念 (Concept) 和实体 边用来表示关系 (Relation) 和属性 (Attribute), 点和边组成知识图谱的基本单位: 三元组 (实体-关系-实体)

到最想要的信息 2) 提供最全面的摘要 3) 让搜索更有深度和广度 应用: 语义搜索 问答系统 推荐系统 多模态知识图谱: 模态: 某种类型的信息或者存储信息的表示形式 传统图谱以文本模态为主,难以描述复杂的现实世界信息 多模态知识图谱可笼统分为属性多模态与实体多模态两大类

属性多模态: 实体: 1) 抽象概念 2) 具有多模态的表示 关系: 实体之间的概念联系 优点: • 易于从现有图谱中进行扩展,即通用 图谱可以 轻松扩展,概念和关系不用改变 • 可以推理视觉知识 缺点: • 实体仍仅限于文本概念描述 • 1 个概念可以对应 N 张图片。 但 1 个图片无法对应 N 张概念

语义信息丰富,场景多源化 • 关系丰富 缺点 • 图谱庞大• 符号复杂

多様杰知识開送的作用:1) 様杰信息搜索 2) 様杰语♡憐竭 応用:推荐系統:引入多様杰信息讲行语♡憐竭

### 知识抽取与表达

<mark>命名实体识别</mark>(NER) 基本概念:识别出文本中的人名、地名等专有名称,和有意义的时间、日期等数量短语等,并加以归类。 命名字体识别是信息抽取中的核心任务,它往往包含两个子任务;判别字体边界 + 判别字体类型 内容: 1) 实体举: 人名 地名、机构名 2) 时间类: 日期、时间 3) 数值类: 货币、百分比 (P.S. 严格定义下哪些不属于命名实体? (部分例子) · 重复指 代的普通名词: 如 飞机、公司 等・人的团体名称以及以人命名的法律、奖项等: 如 共和国、诺贝尔奖 等・ 非时间、日期、货币、 特殊的识别任务:数值识别与检测任务(百分比、钱、邮箱、时间) 时间表达识别与检测任务 难点: 1) 不断有新的命名实体涌现,如新的人名、地名、组织名等 2) 命名实体存在严重歧义 3) 命名实体构成结构复杂,如别名 缩略词、音译、包含数值的实体等 4) 命名实体类型多样:如 John Smith, Mr Smith, John,实际上是共指关系 性能评价:与检索任务大致相同,采用 Precision / Recall / F-value 加以衡量 正确率与召回率的计算方式:方案 1:分子为返回的 正确答案数量 方案 2: 分子为返回的正确答案数量 + 1/3的部分正确答案数量

<mark>命名实体识别方法</mark>:1)<mark>基于词典</mark>:经常作为 NER 问题的基准算法 预先构建一个命名实体词典,出现在词典中的词汇即识别为命 词典的来源:来自于领域公开数据 优点:简单快速,与具体语境/领域无关,容易部署和更新(只需更新词典) 缺点 (与基于词典/匹配分词存在类似问题): ➤ 大部分情况下很难枚举所有的命名实体名 ➤ 构建和维护词典的代价较大 ➤ 週 以有效处理实体歧义 2) 基于规则: I.采用手工构造规则模板,对符合规则的实体进行识别 息、标点符号、关键字、指示词和方向词、位置词 (如尾字)、中心词等 以模式和字符串相匹配为主要手段 可借助结构 该类方法的局限性明显: 不同句式意味着不同的模板 优点: 当提取的规则能较精确地反 具体语言,领域和文本风格,不同领域的句法往往差异极大,如学术卿与二次元、≥ 代价太大,系统建设周期长,移植性差而日需要 基于统计(主流):基于统计模型的分词方法,进一步抽象而言,可以得到一个序列标注问题(四类) 特点: 1. 对特征洗取的要求较高, 需要从文本中 标注: B (词的开始)、M (词的中间)、E (词的结束)、S (单字词)) 选择对 NER 有影响的特征来构建特征向量 2、通常做法是对训练语料所包含的语言信息进行统计和分析,从中挖掘出特征 3、x 语料的依赖也较大,目前缺少通用的大规模语料 4. 大部分技术仍需要进行人工标注训练数据 I.分支一: 基于分类 的命名字体识别方法词性 (part-of-speech) 是词汇基本的语法属性、通常也称为词类。 词性标注就是在绘定句子中判定每个 词的语法范畴,确定其词性并加以标注的过程。词性标注是自然语言处理中一项非常重要的基础性工作。 基本的词性标注方法 (1) 基于规则的方法: 人工或通过大规模语料学习规则(核心思想是按兼类词搭配关系和上下文语填建造词类消歧规则) (2) 基 于统计模型的标注方法: HMM 等面向词序列的方法 (3) 统计方法与规则方法相结合的词性标注方法: 先基于规则排除明显歧义 再基于统计模型标注、最后人工校验 目前、更为简化的是 BIO 标签体系、其中 I 合并了 M、E 所承担的功能、即仅区分开头 和后缀,不区分中间和结尾 单字实体(S)用没有任何 I 后缀的 B 取代 II. 基于序列模型的命名字体识别方法 针对命名字体的类别不同,引入了更多,更细数的标签种类 常 用 模 型 亦 采 用 HMM, CRF以及各种序列深度学习方法(如 LSTM)等 实体对齐/实体匹配:指对于异构数据源知识库中的各个实体,找出属于现实世界中的同一实体。 一般而言,利用实体的属性信息

1.基础任务: 基于表征的知识图谱字体对齐 如何使字体表征拥有统一的向量空间? 联合 学习结构表征和属性字符表征 TransE 结构表征:根据实体关系捕捉两个 KG 之间的实体的相似性 属性字符表征:根据属性 值来捕捉实体的相似性 2.进阶任务:多模态知识图谱的实体对齐 如何表示不同类型的知识:通过多模态知识表征 应用: 文本挖掘, 知识图谱补全, 信息检索和问答系统等

挑战: • 语言的多样性: 同一个实体可能会有多个不同的提及 • 语言的歧义性: 同一个提及可能会对应多个不同的实体 任务框架: 1.提及识别: 其本质与前述分词/实体识别类似 2.候选实体生成 3.候选实体排序 方法: 基于神经网络的实体链接技 术 基于预训练语言模型的实体链接技术 ·<mark>系的表达</mark>:1.二元组 • 适合特定领域关系抽取,类似二部图 2. 三元组 • 适合不定类型关系抽取,最为常见、基础的表达形式

基本关系抽取方法: 1.基于规则的关系抽取 (i.可以根据想抽取的关系的特点 设计针对性的规则,但部分任务可能很难制定规则 ii. 需要领域专家构筑大规模的知识库,代价很大 iii. 知识库构建完成后,对于特定领域的抽取具有较好的准确率,但移植到其他 领域十分困难,效果往往较差) **2. 基于模式的关系抽取** (首先由种子关系生成关系模式,然后基于关系模式抽取新的关系,得到 新关 系后,从中选择可信度高的关系作为新种子,再寻找新的模式和新的关系。如此不断迭代,直到没有新的关系或新的模式产生:

(优缺点: 1.不同算法的差异主要在于模式生成方法和匹配方法。 2.适合某种特定的具体关系的抽取,如校长关系、首都关系。 基于字面的匹配,没有引入更深层次的信息,如词性、句法、语义信息等。 4.移值性差,必须为每一个具体的关系生成自己的识别 模式。) 3. 基于机器学习的关系抽取(采用机器学习方法关系抽取模型,先通过标注语料库训练得到一个判别模型, 再利用该模 往往将关系抽取问题变换为一个分类问题 (二分类或者多分类), 然后采用机 通常采用基于特征或基于核函数的方法加以解决)(开放关系抽取(1)基于知识监督(2)基

基于模式的关系抽取代表方法——DIPRE:其大数思路在于先给定一些已知关系类型的种子实体对,找到出现了这 些实体对的 Occurrences,再学习 Occurrences 的模式(Pattern)。 进而,根据学到的模式,寻找更多符合该模式的数据,并加入到种子集合 中,不断迭代这个过程以实现关系抽取(双重迭代模式关系提取 DIPRE) 基本元素: 1) 元组: 表示关系实例。 如 < Foundation, Isaac Asimov> < Title, Author> 2) 模式:包含常量和变量,例如 ?x , by ?y 的形式 (可表示

"title" by "author" ) 基本假设: 1) 元组往往广泛存在于各个网页源中 2) 元组的各个部分往往在位置上是接近的 3)

基于DIPRE算法,生成模式的基本步骤

共 (Shared) URL前缀与前、后缀。

其他部分采用通配符填充

免因间隔太远而可能导致的语义不相关问题

其次,定义模式如下

urrence归纳为Order (元素的順序) 和Middle (中间部分)

模式: 将同一关系的

・模式的Order和Middle,即为Occurrence集合的Order和Middl

不同实例在网页上所呈现的不同 Occurrence 中,相同内容保留下来,

不同内容采用通配符取代,即可得到近似的模式 URL的前缀 (Prefix)

*所起到的潜在作用*:模式往往仅跟特定网站/体系内,跨网站则模式可能

不适用 被视作从属组织网页: 1) 主机 IP 地址的前三个字段相同,如

182.61.200.X (百度) 2) URL 中的主域名段相同, 如 XXX.ustc.edu.cr

将 URL 的前缀 (Prefix) 引入模式中,用于描述模式的限定范围

模式的Order和Middle。即为Occurrence 集合的Order和Middle

・模式的URLPrefix、Prefix、Suffix,分别为Occurrence集合中最长的公

在表示这些元组时,存在着某种重复的"模式"

- · 首先、输入一组种子元组实例R、如若干-title, suther>的实体对
- 其次,基于种子实例集合R,找到这些元组在阿贡中出现的内容O (Occurrence) , 注意寻找的时候保留上下文信息 (Surrounding Context)
- 进而,基于找到的元组实例O,生成模式P 最后基于生成的模式,找到更多的元组实例。
- 此时可选择停止、戒返回第二步继续基于新实例生成新模式
- . IE PASSONEASTERNORMANICAN

- order = [title,author] (or [author,title])
  a denote as 0 or 1 ddin = "cb> by '

- order = [title, author] (say 0) shared middle = </h>
- shared suffix = (19 pattern = (order, shared prefix, shared middle, shared suffix)
- >>> Foundation <>> by Isaac Asimov (1951)

模式的URLPrefix、Prefix、Suffix、分别为Occurrence集合中最长的公 sho Nightfall sho by Isaac Asimoy (1941) 共URL前缀与前、后缀。

 其他部分采用通配符填充 Snowball: 仅信任支持度(满足每个模式的元组数量)和置信度高(符合该模式的元组确实符合相应关系的概率)的模式

对之间具有某种关系。那么,所有包含这个实体对的句子都是用于描述这种关系,如果某个实体对之间具有某种关系,那么所有包 目的: 获取足够数量的、高质量的标注 语义漂移:不是所有包含该实体 对的句子都表达该关系目错误不断放大 优化方案:1.动态转移矩阵:描述各个类之间相互标错的概率。(随机件较高) 2.规则学习 设计相应否定模式列表去除错误的标签 (对于单一关系的判断,可以通过这种方式进行比较高效的复检。但规则 的生成成本较高)

元程监督:背景:面向文本的关系抽取方法,最大的难点在于获取足够数量的、高质量的标注

6件抽取:事件最信息的一种表现形式。其定义为特定的人、物、在特定时间和特定地点相互作用所产生的客观事实。 事件包含的 要素: 1.事件触发词: 表示事件发生的核心词, 多为动词或名词 (事件触发词的检测与分类是事件抽取的基本任务) 2.事件类型: 与触发 词相对应,往往可以通过触发词分类加以识 3.事件元素:事件的参与者,主要由实体、时间等组成。4.事件元素角色:事件元素在事件 模板: 选定相应模板后, 通过事件元素(事件参与者)与事件元素角色 (事件元素在事件中充当的角色) 的识别, 将相应的元素填入模板合适的槽(描述命名字体基本信息,内容可包括名称/举别/种举)内 **隔定域事件抽取**:与预定义关系抽取存 在相似之外、 即预先定义好目标事件的渗型及每种渗型的具体结构 (包含哪些具体的事件元素) 空切方式・1 平田其干様式 匹配的方法实现 2.采用基于机器学习的方法实现 3.采用远程监督等弱监督的方式,实现限定域的事件抽取 缺陷:种类有限,维护困 开放城事件抽取: 如何在开放域环境下,自动识别未知结构与类型的事件? 思路一: 采用无监督方法 (从而摆脱对于标注语 料的依赖),通过聚类找到潜在的事件簇 思路二:: 与开放关系抽取中的"知识监督"方案类似

图学习任务: 实体/关系的表征学习 关系的抽取与补全 实体/关系的对齐 基于知识图谱的推理 知识图谱应用 基于符号逻辑的知识图谱推理: 基于规则的知识图谱推理 基于符号表示的演绎推理 基于符号表示的归纳推理

### 知识图谱中的类型语言: ABox TBox

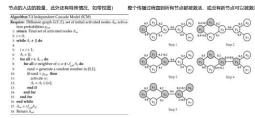
基于表示学习的知识图谱推理: **基于表示学习 (嵌入) 的归纳推理 基于图神经网络的图谱推理 知识图谱上的关系推理任务** (链接预测 实体预测 事实三元组预测) **知识图谱嵌入模型**(1.预测问题与推理评价 • Hit@n:所有预测样本中排名在 n 以内的比例,n 常用的取 求平均 2.TransE: 无法处理一对多,多对一,多对多 3.TransH: 通过向量 wr 将实体投影到关系 r 对应的超平面上,解决一对多 问题 3.TransR: 通过矩阵 Mr 将实体进行坐标转换到关系 r 对应的关系空间中 4. TransD: 关系 r 对应的投影矩阵是动态生成 的,描述同一种关系的不同语义 DistMult RotatE ConvE) 基于图谱推理的推荐技术

知识图谱补全:规则推理、基于表征的方法、基于路径查找的方法、引入大规模预训练模型

<mark>节点的角色</mark>: 1. "意见领袖" 2. "结构洞":在于为组织引入外部的信息(衡量方式: (1) 如果一个节点,移除该节点就会使网络变成 多个连译组件、则该节点则为一个结构词 (2) 聚集系数判定: 某个节点的聚集系数为: 它的任育两个好方也互为好方的概率 (比重) 聚集系数越低,该节点作为中介的作用也就越大 意义:各方沟通的桥梁,结构洞具有排他性)

基本传播模型:1.信息级联发生在有向图上(对于无向图,可以将其连边转化为双向边进行处理) 2. 每个节点仅能将信息传递给与其 直接相连的节点 3. 网络中节点的状态是二元 (Binary) 的 (激活/未激活) (关于激活: (1) 接收到信息, 并且尝试将信息传给别人才叫 激活 (2) 激活能力有一定的时限 (3) 已激活的节点才具备激活其他节点的能力 (4) 激活是不可逆的过程

独立级联模型:每次激活都是一次独立事件,相互不产生影响。同时,每个已激活节点,只有一次机会尝试激活他/他的未激活级民 节点 对于节点 v 而言, 他激活邻居节点 w 的概率采用 Pvw 表示 (基本传播模型里, 为简化考虑, 一般将 Pvw 设为 1/N, N 为 w 节占的入边的数量。此处还有特殊情况。如带权图) 整个传播过程直到所有节占都被激活。或没有新节占可以被激活为止



线性阈值模型:一个用户会被某个信息激活,如果来自他已激活邻居的影响超过某个阈值

阈值预先设定,往往为从[0,1]均匀 分布中随机抽取的一个数值 Algorithm 8.1 Linear Threshold Model (LTM) Require: Graph G(V, E), set of initial activated nodes A<sub>0</sub>

1: return Final set of activated nodes A<sub>∞</sub> while i = 0 or  $(A_{i-1} \neq A_{i+1} \ge 1)$  do

 $A_{i+1} = A_i$ inactive =  $V - A_i$ ; for all  $v \in$  inactive do

if  $\sum_{j \text{ contracted to } v, j \in A_j} w_{j,r} \geq \theta_r$ , then activate  $v_i$  $A_{i+1} = A_{i+1} \cup \{v\}$ 

随机性在干边权重/阈值的确定

Occurrence: 直译为"出现",可以理解成元组在网页中的呈现形式。 and while 一般而言,只有元组的元素在网页中非常接近,才视作 Occurrence,避

线性阈值模型与独立级联模型的区别: 藤机性 1) 对于独立级联模型来说,其随机性在于微硬币的过程 2) 对于线性阈值模型来说,其

<mark>传播最大化问题</mark>: 假定初始的种子节点集合为 S,预期激活的节点集合为 f(S),信息传播最大化的目的,在于限定 S 集合的前提下,最 大化、f(S) 的规模(常见的约束为限定 S 集合的规模、即 |S| 。如果 S 集合中的节点价值不等。则可将约束进一步扩展) 解决传播 最大化问题的启发式思路: 寻找网络中最具影响力的节点 启发式方法: 1. PageRank 及其衍生模型(如 HITS) 2.核心性度量(如度 (Degree)、紧密度 (Closeness)、介数 (Betweenness)) 3.计算单个节点所能够激活的邻居数量。再讲行排序 子梅特性:1. f(S) 函数是非负的 2. f(S) 函数是单调非减的,即 f(S+v)>=f(S) 3. f(S) 函数是具有子模特性 (Submodularity) 的,即: 对于任何 集合对 S. T. 日蓮足 S C T 时、给定节点 v. 有 f(S U (v))-f(S)>=f(T U (v))-f(T) 传播最大化问题是 NP-hard

基于含心質法的传播最大化问题求解伤代码 独立级联模型的松弛版本:稳定状态传播(Steady State Spread, SSS)

Algorithm 7.2 Maximizing the spread of cascades – Greedy algorithm

Require: Diffusion graph G(V,E), budget k

1: return Seed set 5 (set of initially activated nodes) 稳定状态传播模型的核心公式:

 $1 - \pi(i) = \prod_{l} (1 - \pi(l) \cdot p_{li})$ : while i ≠ k do  $v = \arg \max_{v \in V \setminus S} f(S \cup \{v\});$ or equivalently  $\arg \max_{v \in V \setminus S} f(S \cup \{v\}) - f(s)$ 100 表示 / 节点的当前状态、卵板激活的概率 注意此 当 100 = 00寸、星然不影响邻居节点激活状态 处不同 6: S = S ∪ {v};

**衍生传播问题**:独立级联模型的松弛版本:稳定状态传播 节点 状态不再二元化,而是引入一个变量表示当前被激活的概率 如

果被激活概率不为 0,则节点可以持续对外输出信息 / 影响 出现问题: (1) 不合理的反向传播: 解决办法: 1.遍历整个网络 (费 时费力) 2.启发式方法,限制每个节点的迭代次数,超过一定轮次阈值,则该节点被激活的概率 不再更新(阈值大数根据节点到信息源 的最长 / 短路径决定) (2)信息传播与接收的捆绑:解决办法:1.修改目标函数(核心假设:如果某个节点被激活,那么他的所有邻

3.注意力机制: 获取对于整个句子的表示,通过此机制把最能表现关系的句子挑出来(最有效,但依赖于一个高质量的样本集合) 居都被覆盖 贪心算法求解) 2.修改模型框架:单独对信息接受过程进行建模 反爬虫策略:1. User Agent:利用访问网站时浏览器发布的 Request Headers 信息中的 User-Agent 信息,判断用户使用 何种方式浏览。应对策略:利用 Python 的 Request 库允许用户自定义请求头信息的手段。在请求头信息中将 User-Agent 的值改为浏览器的请求头标识, 从而绕开反爬虫机制

IP/账号访问次数/頻率:通过限制特定 IP 地址/账号访问频率和次数进行反爬。其本质在于判断浏览行为是否是人类行为

应对手段: 1)构造 IP 代理池,然后每次访问时蘸机选择代理 (有相关服务,通过各种提供免费 IP 的网站来提供代理池) 2)每次肥取行为后间隔一段时间 3)注册多个账号以保障数据收集 挑战: 账号本身的成本问题、账号被查封的危险 3. <mark>验证码</mark>:通过各类验证码,判断浏览者属于人类还是机器(从简单的字符识别到 复杂的逻辑推理) <mark>应对手段</mark>: 1)简单的字

符识别:基于机器学习与模式识别相关技术 2)复杂的逻辑推理:人工辅助破解 (focus:验证码是一把双刃剑,在抵御爬 4. 动态网页:从网页的 url 加载网页的源代码之后,会在浏览器里执行 JavaScript 程序 (网页内容由脚本加载,而直接抓

取则只能得到空白页面; 此类情况在抓去在线播放的音视频文件时尤为常见) 应对手段:核心思路:模拟调用请求 • 使 用审查元素分析 ajax 请求,如此循环直到获得包含数据信息的 json 文件

 蜜罐技术 • 网页上会故意留下一些人类看不到或者绝对不会点击的链接。由于爬虫会从源代码中获取内容,所以爬虫可 能会访问这样的链接。只要网站发现有 IP 访问这个链接,立刻永久封禁访问者,从而难以继续爬取 <u>应对手段</u>:核心思路 干涉爬虫路径 • 通过工具库判断页面上的隐含元素,使爬虫避开这些元素,可以部分回避蜜罐的诱导。

6.用户权限限制: 不同类型/级别的用户给予不同的内容权限 • VIP、SVIP、蓝钻、红钻、绿钻、各种钻...... • 应对手段: 7. 其他的反爬虫等略。 不同的网页结构:每个相同类型的页面的源代码格式均不相同。 同样也是双刃剑:增加爬取难度

的同时降低用户浏览体验(逼死强迫症) • 多模态的呈现方式: 文字转为图像或视频 • 应对策略: OCR、语音识别、图像

### **网页文字处理** —其干字符匹配的方法(机械分词方法)

正向最大匹配分词 (FMM) 反向最大匹配分词 (RMM) 双向最大匹配分词 (BM: FMM+RMM) 最少切分分词 (最短路 径分词) ASM(d,a,m) d,表示匹配方向,+1为正向,-1为逆向 • a:每次匹配失败后增/减字符数,+1为增字,-1 为减字 • m: 最大/最小匹配表示, +1 为最大匹配, -1 为最小匹配 • 例如, ASM (+, -, +) 即正向减字最大匹配 (即 FMM 方法) 现代汉语最大匹配更实用

FMM: 从左至右尽可能查找最长的词 分词速度较快, 但错误率较高

RMM: 从右至左尽可能查找最长的词 统计证字 RMM 分词效果更好 (由于汉语中偏正结构较多。 若从后向前匹配。可以 活当提高精确度)

BM:综合比较 FMM 与 RMM 两种方法的切分效果,从而选择正确的切分 有助于识别分词中的交叉歧义

最少切分分词方法: 使句子中切出的词数目最少。 • 等价于在有向图中搜索最短路径的问题。 • 将每个字元视作节点, 每 个词形成一条边。 • 边权重可都视为 1, 也可根据词频决定 (尽量切出高频词) • 结合权重/概率之后, 实际上可视作基于 统计的分词方法 N-最短路径法: 保留 N 条最短的路径,以提供更多分词方案 (求一个图中两点的最短 N 个路径: 基于匹配分词方法的优劣 ・ 优点: 效率高、直观性好 ・ 缺点: 对词典的依赖性 ・ 维护高质量词典需要极大的开支 ・ 永远

难以应对新生词汇 • 词汇频率/重要性往往对结果不产生影响

二基于统计的分词方法

# $P(\boldsymbol{w}|\boldsymbol{c}) = \max\{P(\boldsymbol{w}_a|\boldsymbol{c}), P(\boldsymbol{w}_b|\boldsymbol{c}), \dots, P(\boldsymbol{w}_k|\boldsymbol{c})\}\$

一般化过程: 1. 建立统计语言模型 2. 对句子按不同方案进行分词 3. 计算不同分词方案的概率, 选出概率最大的分词结果 理论上、基于统计的分词方法可以不需要词典、但实际应用中第2步可以采用机械分词方法进行分词、以获得候选的分 词集合。• 既发挥匹配分词切分速度快、效率高的特点。• 又利用了无词典分词结合上下文识别生词、自动消除歧义的优

## N-gram 模型

N=1, 一元文法模型 (最大概率模型) , P(w) = P(w<sub>1</sub>) P(w<sub>2</sub>)...P(w<sub>n</sub>) N=2, Bigram模型,  $P(w)=P(w_1)P(w_2|w_1)...P(w_n|w_{n-1})$ 

以Bigram模型为例,基于最大似然估计进行推断

1. 首先,构造训练语料库,计算所有的 $C(w_n)$ 与 $C(w_{n-1}w_n)$ 。 2. 其次,对于每一个可能的分词序列w,计算以下公式

**Bigram:**  $P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$ 

・ 其中, $\mathcal{C}(w_{n-1}w_n)$ 指词序列 $w_{n-1}w_n$ 在语料库中出现的次数。 ・ 其中, $P(w_n|w_{n-1}) = \frac{\mathcal{C}(w_{n-1}w_n)}{\mathcal{C}(w_{n-1})}$ 

而C(w<sub>w</sub>)指某个单词w<sub>w</sub>在语料库中出现的次数。

3. 最后,返回最大的P(w)所对应的分词序列作为结果。

当N=1时, N-gram模型退化为一元文法模型, 此时词与词之间是独立的。

•  $P(w) = P(w_1) P(w_2|w_1)...P(w_n|w_{n-1})$ 

$$P(W) = P(w_1, w_2, ..., w_i) \approx P(w_i) \times P(w_2) \times ... \times P(w_i)$$
   
    
**独立性假设**,一元文法

 $P(w_i) = \frac{w_i$ 在语料库中的出现次数n 语料库中的总词数N

基于统计文法模型的优劣: • 优点:减轻了对于词典的依赖性 • 然而,这种依赖并非完全消除,取决于性能与效率的平衡 如果深度结合机械分词(匹配分词),则效率提升但依赖词典
 如果减少对词典的依赖,则需要更多地遍历潜在的组合(新 空间巨大!) • 缺点:依赖已有数据中词频的统计,对于新生词汇或专业词汇不友好 • 冷门领域的稀有词汇往往难以准确

# 划分 · 易受数据集先验偏差 (Bias) 的影响

隐马尔可夫模型的五个核心要素:两个集合、三个矩阵

三个矩阵: · 初始状态概率矩阵: 第一个字属于某种隐含状态 (BMES) 的概率 · 隐含状态转移概率矩阵: 各种隐含状态 (各种标签) 之间的转移概率 • 观测状态概率矩阵: 从隐含状态 (标签) 到观测值 (字符) 的转移概率

<mark>專用词</mark>: Stopwords,指文档中频繁出现或对实际语义影响不大的词语。去停用词原因:• 重复率很高,会造成索引中的倒排表很长 影响查询性能。 • 对最后结果的排序没什么贡献,反而可能产生干扰。

常用的停用词识别方法。 较为成熟的停用词识别方法有: 文本概率、词频统计、统计算等。。 更为复杂的算法将结合统计与句法或内 现代搜索引擎的趋势是逐渐减少对停用词的使用。 • 现代搜索引擎更关注利用语言的统计特性来处理常见词问题。 • 采用压缩技术

降低停用词表的存储开支。 • 引入词项权重,将高频词的影响降至最低。 • 索引去除技术,低于权重的词项将被排除 O

<del>归一化处理</del>: 归一化/词根化,指还原词语的特殊形式的过程。可以有效降低词项的数量并减少歧义

词干提取: Stemming, 指去除单词前后缀, 获得词根的过程。常见的前后词缀有"复数形式"、"过去分词"、"进行时"等 词形还原: Lemmatisation, 指基于词典, 将单词的复杂形态转变成最基础的形态。并不是简单地将前后缀去掉, 而是会根据词典将单

词干提取与词形还原的相同点: 1)目标一致。两者的目标均为将词的不同形态简化或归并为基础形式。2)结果交叠。两者不是互斥关系

•词干楼取与词形环原的不同点 1)在原理上,词干楼取采用"统诚",而词形环原采用"转夺"。 2)在复杂性上,词形环原需考虑词

缀转化、词性识别等,更为复杂。 3)在实现上,词干提取主要利用规则变化,而词形还原更依赖于词典。 4)在结果上,词干提取不-基本词干提取方法: Porter Stemming 基本流程: 1. 去除单词的复数形式 2. 去除 -ed (ly) 或 -ing (ly) 等后缀 3. 将 -y 改为 -

4. 处理双重后缀、如 -ization 等 5. 处理 -full , - ness 等后缀 6. 处理 -ant , - ence 等后缀 7. 去除掉最后的 - e 和 - II

### 倒排表 (Posting List): 文档 ID 列表,列举词项在哪些文档中出现 学録表引 (kov) Document Term 化占 思维护 松建 维护成大家 缺点 搜索鲜时长 搜索快

倒排表三步走: 1.检索每篇文档, 获得对, 并写入临时索引 2. 对临时索引中的词项进行排序 3. 遍历临时索引, 对于相同词项的文档 ID 讲行合并 基于倒排表的查询:如果两个倒排表的长度分别为 x 和 y,则合并供需 O(x+y)次操作 同时扫描两个倒排 表,所需时间与倒排记录的数量呈线性关系 优化问题: 1. AND 是合并操作,先处理文档频率小的,再处理大的 2.任意组合的布尔查询: 首先,获得所有词项的文档频率 • 其次

保守地估计出每个 OR 操作后的结果大小 • 考虑 x+y 的最坏情况 • 最后,按照结果从小到大的顺序执行 AND 3.跳表指针:1)如何诊 ■ I.较多的指针+较短的步长→更多的跳跃机会(耗费更多的存储空间) II.较少的指针+较长的步长→更少的指针比较次数(存储 空间消耗更少,但跳跃机会也更少) 2)如果倒排表长度为 L,则间隔 L^(½)均匀放置跳表指针 3) 索引的动态变化也会影响跳表指针的

设置 • 如果索引相对固定,建立有效的跳表指针就相对容易 • 反之,经常更新的索引很难建立合适的跳表指针

倒排表还可能面临更多需求 (词组查询、词项邻近、文档区域)

汇表迅速增长): 将文档中每个连续词对看成一个短语,构建面向二元词的倒排表,并处理两个词(或多个词)构成的短语查询 更长的 短语查询:可以分成多个短查询来处理(如果采用二元词索引拼接的方式,对于该布尔查询返回的文档,不能确定其中是否真正包含最 2) 位置信息索引:在记录词项的同时,记录它们在文档中出现的位置,可以达成更广泛的查询。对于 短语查询,仍采用合并算法(AND),查找符合的文档。当然不只是简单地判断两个词是否出现在同一文档中,还需要检查他们出现的 位置情况是否符合要求。并且该方法能够用于邻近搜索(例如,间隔 k 个词)

1)除了文档 ID,搜索引擎往往在倒排表中加入更多元素(词项频率(Term Frequency)、词项类型等) 2)除了基本的词查询之外

短语查询的需求解决方案: 1) 二元词索引 (问题: 证

索引存储: 词典与倒排表分开存储 (可以支持并行、分布式查询) 词汇表存储: 顺序、哈希、B/B+树、Trie 树 (前缀树)

(节省磁盘空间(¥),提高效率(内存利用率或数据传输速度)前提:快速的解压缩算法。目前的启发式算法效率都比较高。)意义: 1) 对词典: 压缩得足够小, 可以直接放入内存中, 提升效率 2) 对倒排记录表: 减少所需的额磁盘空间, 可以更多移入内存, 并且可 词典基本存储方式: 定长存储(词项本身设置定长为20字节,另外需要记录文档频率和指向倒排表各4字节),但是空间会浪费(英文

压缩词项列表: 将词典视作单一字符串(在词项字符串中,每个词项平均长度8个字节•其次,词项文档频率与倒排表指针各4字节 不变 · 最后, 词项指针约 3 字节 按块存储: 单一字符串在词项指针上需要占用较多额外空间。通过为每 k 个词项存储—个指针,来减少指针的总数量。需要额外 1 个字

节用于表示词项长度。例如,当块的大小 k=4 时 • 不采用按块存储时,每个指针花费 3 字节,共需 12 字节 • 采用按块存储时,只需 要花费 3+4\*1=7 字节 在搜索上的问题:二分查找只能在块外进行,块内采用线性查找方式(随着 k 上升,线性查找部分增多 效率更低) 改进: 前端编码 (Front Coding) 采用间距代替文档 ID----再进一步: 可变长度编码 (第一个存原值, 其它存间距!

先存储G. 并分配1bit作为延续位 →8automat\*a10e20ic30ion 如果G<128、则采用第一位延续位为1+7位有效二进制编码的格式。 如果G>=128、则先对低阶的7位编码、然后采取相同算法对高阶位进行 编码。最后一个字节 (8bit) 的延续位为1, 其他字节延续位为0. 214577 的二讲制为 1101/0001100/0110001.

VB 编码为 00001101/00001100/10110001

面向单查询的基本评价指标: 1) 准确率 (Precision): 指检索出的文档中 • P/N: Positive or Negative, 表示算法对样本的实际 ・ TF: True or False,表示算法判断的正确与语(Ridi未是否一数)相关文档所占的比例,也称查准率。计算公式为TP/(TP+FP) 2)<mark>召回率</mark>

TP: Time or False。表示專品判斷的正确与自信(相應果然 但种國事的含义: - IT: Time Novation,样本为正例,目睹對從为正,將真正 - IN: False Negative,样本为正例,而國政總對原之份,將假 - IT: False Negative,样本为正例,但國政總對原立为此,即居 - IN: Time Negative,样本为条例,且被判除力灸,即真灸

(Recall):指所有相关文档中,被检索出来的部分的比例,也称查全率。计 算公式为 TP/(TP+FN) 召回率的近似计算: 缓冲池 (Pooling 针对某一检索问题,各个算法分别给出检索结果中的 Top N 个文档,将这些 结果汇集起来并进行人工标注,从而得到一个相关的文档池(虽然它实际上仍

然无法得到全部相关文档,因此并不能得到召回率的绝对值。但是,它可以比较各个算法的相对优。N 通常取 50-200)

sec),即准确率与召回率的加权调和平均数 P-R 曲线:可以直观地看出准确率 与召回率之间的平衡关系  $F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$ ROC 曲线(Receiver Operating Characteristic Curve,接受者操作特征曲线)以真正

通常情况下、我们取(-0.5或2-1 (即两者同等重要) 此时,可得基本的F1值,即F=2PR(P+R)

率【TP/(TP+FN)】和假正率【FP/(FP+TN)】作为两条轴线。对角线表示区分能力为 0, 即随机猜测。在对角线上端越远,效果越好。低于对角线的结果无意义 (无区分度)

· 如果线A将线B完全包住,显然线A对应的算法效果更好 P-R曲线与ROC曲线的选择 • ROC曲线兼顾正负样例,更为全面,而P-R曲线则只考虑正例 如果两条线发生重合,则可依据以下规则判别: · 计算AUC, AUC更高者效果更好 用户往往更关心正样本,如果面向特定应用场景(如检索),P-R曲线是个好选约 单份数据正负样本比例失调时,P-R曲线更合适 · Area under curve, 即曲线下面积

 可通过积分近似计算 当色样本比重过高时、色例的数目众多致使FPR的器长 不明显,导致ROC曲线呈现一个过分乐观的效果估计 从而难以体现出性能的差异性 • 另外,当使用P-R曲线时,可使用平衡点计算 · 平衡点即Precision - Recall的点。值越高越好

P-R曲线与ROC曲线的选择 · P-R曲线受分布影响大,多份数据且正负比例不一时ROC曲线更合适

• ROC曲线两个指标各自针对正负样本,而Precision只针对正样本,受影响较大

P@N,即 Precision@N,指前 N 个检索结果文档的准确率:如果相关文档数小于 N,P@N 的理论上限必定小于 1 R@N,即 Recall@N,指前 N 个检索结果找回的相关文档比例 : 1) 由于返回结果有限,Recall@N 值,甚至其理论上限往往都远

小于 1 2) 理论上限为 N/相关文档数,即使通过 Pooling 加以控制仍然较小

R-Precision=P@相关文档总数 未插值AP: 查询0共有6个相关结果,排序返回了5篇相关文档,位置分别是第1 第2,第5,第10,第20位,则AP~(1/1+2/2+3/5+410+5/20+9)/6 (注意补o)

插值AP: 事先选定插值点数并进行插值。例如,当我们计算11点平均时,计算在召回率分别为0(第一条)。10%。20%.....100%的十一个点上的正确率求平均 第化AP: 只对返回的相关文档进行计算。AP=(17+2/2+35+4/10+5/20)5、倾向 那些快速返回结果的系统,没有考虑召回幸和补零的情况 (不补o)

累计增益 CG:位置 1 到 p 的检索结果相关度之和;  $DCG_p = rel_1 + \sum_{i=2}^p rac{rel_i}{log_2(i)} \ ext{(ppt用第一个)}$  或  $DCG_p = \sum_{i=1}^{p} \frac{2^{rel_i} - 1}{\log_2(i+1)}$ 

序列,并对此序列计算 DCG 值所得到的。 理想的输出结果序列为: 3, 3, 3, 2, 2, 2, 1, 0, 0, 0 - 曲此计算:DCG依次为: 3, 6, 7.89, 8.89, 9.75, 10.52, 10.88, 10.88, 10.88, 10.88 而与此同时,基本的DCG结果如下:

NDCG = DCG / iDCG

· 3, 5, 6.89, 6.89, 6.89, 7.28, 7.99, 8.66, 9.61, 9.61 (単調非波特性) 由此可得 NDCG结果如下: 1, 0.83, 0.87, 0.76, 0.71, 0.69, 0.73, 0.8, 0.88 · 可以看到任何查询结果位置的NDCG值都规范化为[0,1]之间的值

最大边界相关性 (Maximal Marginal Relevance, MMR)

多查询指标: 1) MAP: 对所有查询的(简易)AP 求算数平均(可以反映全部查询的综合效果,但在查询难度不平衡的条件下有 误导) 2) MRR: 多个查询所得的倒数(第一相关文档)排序求评价(eg. 两个查询,第一个的第一个相关文档在位置 2,第二个的第一 个相关查询在位置 4,则 MRR=(1/2+1/4)/2=3/8, 即平均在 8/3 的位置上找到第一个相关文档) 3) 方差: 一个检索系统对不同查

询的方差通常大于多个检索系统对相同查询的方差,说明查询难度差异大 多样性评估:基本形式:给定一个查询 q,返回一个多样化的结果文档集合 R(q)(R(q)作为一个整体,应满足以下条件: • R(q)中 所有的结果文档都与查询 q 本身有较大的相关性。 • 总体上要有较小的冗余度,以覆盖 q 的不同方面。) 核心思想:降低用户

无法获得所需信息的风险(尽可能确保排序靠前的结果中至少有一个结果满足用户的需求) 总体需求: 衡量不同文档之间的主 题差异性。两种衡量方式: 1) 隐式模型: 只计算文档之间的差异性・文档是什么内容。不会也无法进行详细考量 2) 显式模型 更加具体地考量文档所对应的用户意图 • 会从文档中抽取主题,并显式地实现主题的多样化 · 显式模型代表: FM-LDA

- 最早的相关性衡量模型,由Carbonell和Goldstein在1998年提出 • 由Carterette与Chandar在2009年提出,考虑文档的不同子主题(Facet)  $MMR^{def} = Arymax_{d_i \in RX} |\lambda P(d_i \mid q) - (1 - \lambda) max_{d_i \in R} P(d_i \mid d_f)$  $L(y_j(F, D) = \prod_{i=1}^{n} \left(1 - \prod_{i=1}^{n} (1 - p(f_j \in d_i))^{\gamma_j}\right)$ · 该公式由两部分组成。其中: · 前一半表示文档集与查询 q 的相似性

 其中、p(f, s.d.)表示文档Di包含主题Fi的概率、Yi=1表示某文档被选中 同一半表示文档之间的多样件 由此,联系部分表示主题与至少被一个文档所涵盖的概率 λ用于資节两部分之间的比例 · 文档可以采用(f-idf, 或者word2vec等文本表征工具进行表征 相关性反馈:用户在查询后标记相关/不相关,然后迭代更新查询,以获得更好的结果 过程:1. 首先,用户提出一个查询条件

· Facet model with LDA (FM-LDA)

(Query) 2. 对于返回的文档,用户标出相关与不相关的部分 3. 系统根据用户反馈,获得用户信息需求更为准确的描述 (a) 复 于相关性反馈。更新查询条件 b) 基于新查询条件。获取新的线果文档并再次提交用户进行评估) ➤ 上述过程将根据情况进行一次 目的: 通过相关性反馈, 获得用于表达用户查询意图的最优查询条件 (常见方 式: 为已有词项 添加不同权重,或增加 新的词项 · 这一过程应对用户隐藏) 相关性反馈类型: 1) 显式反馈: 用户显式地 参与交互过程(如用户点击记录,但是只有正样本,并且用户不点击,不代表完全不相关。所以有了拓展的显式反馈: 收集负面评 价的渠道日益丰富,还有用户评论) 2) <mark>隐式反馈(更常见)</mark>:系统追踪用户行为来推测返回文档的相关性(判定不一定很准确, 但省却了用户的显式参与行为。如鼠标键盘动作,用户眼球动作等)优点: • 不需要用户显式参与,减轻用户负担,提升用户体验 • 用户行为某种程度上可以反应其兴趣,因此具有可行性 缺点: • 对行为分析有着较高的要求 • 准确度难以保证 • 某些情况下 需要增加额外设备(且很贵!) 3) 伪反馈:在没有用户参与的前提下,直接假设返回结果是相关的,并进行反馈(准确性差,可 能查询漂移) 4) 查询扩展: 用户针对词项的合适程度给出反馈,这些反馈将被用来构建更为完整的查询条件。 • 暗含的功能: 用 户选择和确认的音询扩展能够更好表达其音询意图。 字现: 利用同义词辞典, 可以字现音询条件的扩展。一般而言, 音询 扩展有助于提升查询的召回率(找得更全)。但是,可能会影响准确率,尤其在扩展词存在歧义的情况下。编纂和维护同义词词典常 要很大的代价 。 类型: 1.利用人工编纂的同义词词典 2.全局分析与同义词词典的自动生成 3. 基于搜索日志进行 自动构建相关词词典的两种思路(通过分析文档集中的词项分布,来自动生成同义词/近义词词典): 1.如果两个词经 常和相似的词共同出现,则它们很可能是相似的 2.如果两个词经常与相同的词以一种特定的语义关系共同出现,那 么他们很可能 具相似的 基于搜索日志的音询扩展: 搜索日志目前是搜索引擎音询扩展的重要方式 (潜在扩展音询 1 音了一个词又音写-个 (后是前的) 2.查了两个词 (互为) 后都点进同一个链接) 相关词扩展的潜在问题: 词项关联的质量是一个问题 (1. 有歧

## 排序的难点: Web 网页的信息组织与内容质量参差不齐,用户庞杂且缺乏知识和经验,用户意图多样差异巨大

够 获得更多的相关文档)

好的排序: 网页内容匹配程度+网页内容质量 信息检索模型:文档与查询的表示形式与相关性的框架。实质是对文档基于相关性进行排序,在理解用户的基础之上产生近似用户

决策的结果从而在顶部返回最相关的信息。[D,Q,F,R(Di,q)] 布尔模型 D:文档表达-词项的组合; Q:查询表达-布尔表达式; F:完全 (二值)匹配: R.满足布尔表达式,相关性为 1(否则 0

JACCARD(A, B) =A U B 词项频率 TF(t,d): 词项 t 在文档 d 中出现的次数 相关度计算: Jaccard 系数: 对数词频: $wf_{t,d} = \begin{cases} 1 + \log_{10} t f_{t,d} \end{cases}$  $tf_{t,d} > 0$ otherwise

$$W_{t,d} = \left(1 + \log t f_{t,d}\right) \cdot \log rac{N}{df_t}$$
 成为  $rac{1}{1 + \log t}$  公司  $\frac{N}{2}$  公司  $\frac{N}{$ 

文档表达,每个文档可视作一个向量,其中每一维对应词项的 tf-idf 值;Q:查询表达,可视作一个向量,其中每一维对应词项的 tf-idf 值; F:非完全匹配方式; R:使用两个向量之间的相似度来度量文档与查询之间的相关性) 优点: 简洁直观, 可支持多种不同度 量和权重方式,实用效果不错; 缺点:缺乏语义层面的理解和匹配,同时依赖 tf-idf 值也可能造成干扰,用户无法描述词 项之间的关系,词项之间的独立性假设实际上不成立。

文档匹配与迭代优化

欧氏距离 (不行,太逊了。对于向量长度 (文档长度) 非常敏感  $dist(\vec{x},\vec{y}) = \sqrt{(x_1-y_1)^2 + (x_2-y_2)^2 + \ldots + (x_n-y_n)^2}$ 非常被害  $dist(x,y,r) = \frac{\sum_{i=1}^{pr} q_i d_i}{\sum_{i=1}^{pr} q_i^2} \frac{1}{\sqrt{\sum_{i=1}^{pr} q_i^2}} \frac{1}{\sqrt{2}} \frac{1}{\hat{q}_n} = \alpha \hat{q}_n + \beta \frac{1}{|D_i|} \sum_{j_i \neq j_i} \hat{d}_j - \gamma \frac{1}{|D_{ir}|} \sum_{j_i \neq j_i} \hat{d}_j$ : 余弦相似度 (文档向量与查询向量的夹角大小来计算): Rocchio (罗基奥) 算法 (使查询尽可能离与之相关的文档更近,离与之不相关的文档更远)

基中、Dr.为已知相关文档的向量集合、Dnr.为已知不相关文档的向量集合 • α0 为初始音询向量。α β、ν 为权重、根据手工调节或经

基于迭代的音询意图更新 : 正反馈 vs 负反馈

• 正反馈的价值往往大于负反馈 用户更关心符合需求的标准答案,而不是错 Original query 0 4 0 8 0 0 α = 1.0 0 4 0 8 0 0 误答案。相应的,可以通过设置 β > γ 来给予正反馈更大的权重 Positive Feedback 2 4 8 0 0 2  $\beta$  = 0.5 1 2 4 0 0 1 (+)

 很多系統甚至只允许正反馈,即 γ = 0 · 收集真正的负反馈往往比较困难 Negative feedback 8 0 4 4 0 to y = 0.25 2 0 1 1 0 4 1

Word2vec: 1. 根据上下文预测中心词 (CBOW) 2. 根据中心词预测上下文 (Skip-gram) CBOW算法详细步骤:

New query -1 6 3 7 0 -3

基于深度学习的CBOW算法详细步骤: 使用one-hot (0/1) 向量表示词項: x<sub>i</sub>

计算的数据(深度を示する場):

词, 复杂度约为 O(KV), 即考虑窗口

计算交叉熵损失,反向传播优化模型

两个模型比较: 1) 从性能上说 • CBOW 模型仅预测中心词,复杂度约为 O(V),即词表规模 • 而 Skip-gram 模型基于中心词预测周边

 $L = -\sum_{v} v_v \log(\widehat{v}_v)$  基于多个不同量动物口进行多次训练 最终得到各个证证的表征: Wy

从效果上说 • 在 Skip-gram 模型中,由于每个词都可以作为中心,都将得到针对性训

计算预测的中心团的概率

练 · 因此,对于生僻词(数据稀疏)的训练而言, Skip-gram 模型效果更好 Word2vec 优缺点: 1) 优点 • 有效表征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 2) 缺点 · 无 法解决一次多义的问题,例如 play music 和 play football • Word2vec 是一种静态的方式,其词项表征一旦训练确定就不会再做更

排序学习算法: Pointwise 算法 Pairwise 算法 Listwise 算法 Pointwise: 基本假设: 训练样本中的任何一个查询-文档对,都可以映射到一个分值或一个有序的类别(如优良中差) • 相应的,给定

改。 因此,虽然通用性强,但是无法针对特定任务进行动态优化

一个查询-文档对, Pointwise LTR 将试图预测其分值/类别。 · 常见的模型类别包括: 1) 回归, 将查询-文档对映射到具体分数 2) 分类,将排序问题转化为一个面向有序类别的二分类/多分类问题 3) 有序回归,在映射到具体分数的同时保持样本之间的有序关系 局限性: 1) 首先, Pointwise 类方法往往更为注重文档的相关度得分, 而并 不注重文档之间的相关性排序(Pairwise 类方法的出现,为解决这一问题提供了新的手段)2)其次,不同查询所对应的文档,尤其相关

Pairwise: 基本假设:同样是将排序问题转化为分类问题(二分类或三分类) • 每次比较一个查询与两个文档,衡量两个文档的偏序 (Partial Order) • 分类器的目的在于判断哪个文档应该排在前面(对应的标签为{1,-1}或{1,0,-1}(0表示两个文档可以并列)) 相比于 Pointwise 类算法,Pairwise 类排序算法通过衡量样本之间的偏序关系,实现了从绝对相关性(分值)到相对偏序的改进 缺陷: 1) 两两成对导致样本数大为提升, 计算资源开支增加 2) Pairwise 类算法仍然受样本不平衡问题的影响 3) 最后, Pairwise

类算法无法体现全局排序的合理性。综上,由此引出了最后一类算法: Listwise 类算法 Pairwise 局限性: 样本不平衡性的影响以及 Listwise:基本眼路:直接面向整体排序结果进行优化(直接将排序的完整队列作为学习的对象)解决这一问题采用以下两种思路:1)

直接采用某种 IR 指标对排序进行优化 2) 直接设计面向完整排序的损失函数 代表性思路:采用某种 IR 指标对排序进行优化 通常情况下,由于全局优化的作用,Listwise 类排序算法可以取得相比于 Pointwise 和 Pairwise 类算法更好的效果。 小的挑战:

例如两个网页并列的情况(相当于 Pairwise 中文档对标签为 0 的情况) 拓展: Listwise 与 Pairwise 的融合: FocusedRank

 $PR(p_j)$  PR(pi)为网页 pi 的 PageRank 值,初值为 1/N; PR(pi)为指向网页 pi 的某  $d\sum_{p_j \in M(p_i)} \frac{1}{2}$  $PR(p_i) = \frac{1-\alpha}{N} + d$ 个网页 pj 的 PageRank 值; L(pj)为网页 pj 发出的链接数量; d 为阻尼系数,  $L(p_j)$ 取信在 0-1 之间: N 为网页总数, M(pi)为链入 pi 的页面集合 邻接铂阵->跳转铂阵·转置(行列互换)后各数值除以出边数量

则有Pn+1 = APn, 形成马尔科夫过程





对于某个Topic来说,起点仅限于隶属于该Topic的网站 (I-d)e/N更改为了(I-d)s/|s|, s为一个N维向量 · 如果某个网站隶属于该主题,则 s 中的该维为1,反之为0 ・ |s|表示 s 中为 | 的元素的个数。 孤立结点:没有任何入边,单纯使 陷阱节点: 只有一条指向自己的边,没有其它出边; 终止节点:没有出边,如同黑洞;

用 markov 链无法跳转,仅有初始概率,不再更新也不影响其它结点。 解决: 重启动: d 或 alpha (随机跳转概率是 1-d)。 why?其中的(1-d)/N 的部分相当于以一定等概率被选中作为新起点,由此跳出了陷阱和黑洞的干扰(所有节点全联通)。 马尔可夫过程的三个收敛条件(转移矩阵 A 所有元素非负,且每一列元素和为 1,为马尔科夫矩阵;Restart 保障圈是强连通,则 A 不可  ${\mathbb R} A = dM + [(1-d)/N]ee^T$ ,其中M为跳转矩阵,c为所有元素都为I的列向量 解特征方程(A-I)x=0 且 $\sum x_i = 1$  得概率 $x_*$ ( $\lambda=1$ )