

事理图谱: 描述逻辑社会, 研究对象是谓词性事件及其内外联系 (区别), 借助事理逻辑链接形成对于事件的推理 (应用)。不足: 关系的边界比较难以界定; 对于事件关系的研究大量集中于因果关系等, 对于其他关系的研究较少; 事件的定义不够明确, 从而引起抽取事件比较困难。
多模态图谱: 实体和属性可能是多模态的 (区别), 表示与整合多模态知识 (应用)。作用: 模态知识互补: 不同模态协同服务于实体的理解; 模态信息搜索: 以不同模态作为信息入口; 模态语义增强: 利用知识图谱增强多模态任务; **属性多模态**: 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: 实体不仅限于文本概念描述, 1个概念可以对应N张图片, 但1张图片无法对应N张概念。**实体多模态**: 优点: 视觉语义信息丰富, 场景多源化; 关系丰富。缺点: 图谱庞大; 符号复杂。

Lec 12 知识抽取与表达 (上)

命名实体识别是信息抽取中的核心任务, 它往往包含两个子任务: 判别实体边界; 判别实体类型。难点: 与分词的难点非常相似: 不断有新的命名实体涌现; 命名实体存在严重歧义; 命名实体构成结构复杂; 命名实体类型多样 (如共指关系复杂)。
命名实体识别方法: **基于词典**: 预先构建一个命名实体词典, 词典中的词汇即识别为命名实体。优点: 简单快速, 与具体语境无关, 容易部署和更新。缺点: 难枚举所有的命名实体名; 维护代价, 实体歧义。
基于规则: 手工构造规则模板, 对符合规则的实体进行识别。优点: 当提取的规则能较精确地反映语言现象时, 性能较好。缺点: 不同表达对应不同规则, 规则库大; 规则往往依赖于具体语言、领域和文本风格, 代价太大, 系统建设周期长、移植性差。
基于统计: 抽象为序列标注问题。四类标注: B (句的开始)、M (句的中间)、E (句的结束)、S (单字词)。分支一: 基于分类的命名实体识别方法, 将NER视作一个多分类问题, 通过设计特征训练分类器的方法加以解决。分支二: 基于序列模型的命名实体识别方法, 与分词中的序列标注方法思路类似。区别在于标注的不同。
实体对齐: 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 (如路社交网络用户匹配)。基于表征的知识图谱实体对齐: 利用相似性合并使关系表征拥有统一的向量空间。
实体链接: 将文本中的提及链接到知识库中的实体上。方法: 神经网络、预训练语言模型。

Lec 13 知识抽取与表达 (下)

知识抽取: 从文本中识别出两个实体 (或多个实体) 之间存在的事实上的关系。意义: 搜索引擎发现和关联知识的重要渠道, 知识库构建与知识关联的基础性手段, 是支持问答系统、推荐系统等应用的有力工具。方式: 1. 基于规则; 2. 基于模式; 3. 基于机器学习。
基于规则: 优缺点: 1. 通常针对特定领域的特定关系抽取任务, 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; 2. 基于手工规则的方法需要领域专家构建大规模的知识库, 这不但需要有专业技能的专家, 也需要付出大量劳力, 因此这种方法代价很大。3. 知识库构建完成后, 对于特定领域的抽取具有较好的准确性, 但移植到其他领域十分困难, 效果往往较差。
DIPIRE: **基本元素**: 元组: 表示关系实例, 如-**Foundation**, Isaac Asimov>-<-**Title**, Author>-。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: 元组往往广泛存在于各个网页面上; 元组的各个部分往往在位置上接近的, 在表示这些元组时, 存在着某种重复的模式。**流程**: 首先, 输入一组种子元组实例 R, 如若干<title, author>的实例对; 其次, 基于种子实例集合 R, 找到这些元组在网页中出现的內容 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 R, 此时可选择停止, 或返回第二步继续基于新实例生成新模式 (此时生成的新模式可能与之前的模式有所差异)。
Occurrence: 元组在网页中的呈现形式。一般而言, 只有元组的元素在网页中非常接近。**模式**: 将同一关系的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。将 URL 的前缀 (Prefix) 引入模式中, 用于描述模式的限定范围。**生成模式的三步骤**: 首先, 将 Occurrence 归纳为 Order (元素的顺序) 和 Middle (中间部分); 其次, 定义模式如下: 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 集合中最长的公共 (Shared) URL 前缀与后、后缀, 其他部分采用通配符填充。
远端监督: 思想: 如果某个实体对之间具有某种关系, 那么, 所有包含这个实体的句子都是用于描述该关系。**局限性**: 语义漂移, 不是所有包含该实体对的句子都表达该关系, 错误模板会导致关系判断错误, 并通过不断迭代放大错误。**优化方案**: **动态转移矩阵**: 引入一个动态转移矩阵, 描述各类之间相互标错的概率。在利用算法得到的关系分布的基础上乘以这一转移矩阵, 即可得到相对更为准确的关系分类结果。**规则学习**: 模拟远端监督的启发式标签过程, 设计相应的否定模式列表 NegPat(n), 专门用于去除错误的标签, 即某些关系的判断是否为错误。注意力机制: 即使是被打入同一个包里的句子, 不同句子对于训练关系判别模型的贡献度也不相同, 这一贡献度可以采用注意力模型加以衡量。采用深度学习技术, 获取对于整个句子的表示。进而, 通过注意力机制, 将最能表达这种关系的句子们挑选出来。
事件抽取: 事件是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。**基本要素**, 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别; 事件元素: 事件的参与者, 主要由实体、时间等组成; 事件元素角色: 事件元素在事件中充当的角色。

Lec 14 知识图谱与图计算

图表示学习算法: 将图数据进行向量化表征, 映射到一个低维的向量空间, 在这个低维向量空间中, 图的结构特征和语义特征得到最大限度的保留。**邻接矩阵**: 每一行表示一个节点, 1/0 分别表示与对应节点是/否连接, 这一行可视为该节点的一个表示向量; 该思路可用于最基础的图聚类问题; 局限性: 未能充分融入节点结构信息, 节点属性信息无法加入。**基于随机游走的图表示学习**: 基于随机游走的邻居节点序列, 挖掘图结构信息。**基于图神经网络的图表示学习**: 利用神经网络来学习图结构数据, 提取和挖掘图结构中的特征和模式。图神经网络技术: 路径游走、面向表征整合的消息传递。

Lec 16 社会网络

独立级联模型: “独立”体现在, 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/她的未激活邻居节点。如果某个节点在第 t 轮被激活, 那么, 他仅有一次机会, 即仅能在 t+1 轮, 尝试激活他所有未被激活的邻居节点。t=1 时, 仅有种子节点可以尝试激活其他节点。整个传播过程直到所有节点都被激活, 或没有新节点可以被激活为止。
对于节点 v 而言, 他激活邻居节点 w 的概率采用 P_{vww} 表示。基本传播模型里, 为简化考虑, 一般将 P_{vww} 设为 1/N, N 为 w 节点的入边的数量。
线性阈值模型: 将信息传递过程视作多人影响的叠加过程。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 (或根据用户对信息的兴趣等决定)。
对比: 线性阈值模型与独立级联模型的区别: 随机性。对于独立级联模型来说, 其随机性在于抛硬币的过程。独立级联模型是完全随机过程, 每一次的结果可能都不相同, 一般需要重复多次以确定个体节点被激活的可能性。对于线性阈值模型来说, 其随机性在于边权重/阈值的确定。如果采用启发式方法确定边权/阈值, 则该方法结果完全由方法设计决定, 一旦确定边权/阈值 (无论何种方式), 其结果具有唯一性。

信息传播最大化方法: PageRank 及其衍生模型、核心性 (Centrality) 度量、计算单个节点所能够激活的邻居数量, 再进行排序。存在问题: 在寻找“最具影响力的节点”时可行, 在确定影响力节点集合时不可行 (可能存在影响范围重叠)。

目标: 找到一节点集合 S, 使得 f(S) 的期望最大, 并且 |S|=k 只选择 k 个节点作为初始节点。f(S) 特性: 子模特性。1.f(S)非负; 2.f(S)单调非减, f(S+v) >= f(S), 新增加一个节点, 至多不增加新激活, 不至于减少; 3.f(S)具有子模特性: 对于任何集合对 S,T 且满足 S ⊆ T 时, 给定节点 v, 有 f(S+v)-f(S) >= f(T+v)-f(T)。

在 ICM / LTM 等模型定义下传播最大化问题可以归约为集合覆盖和节点覆盖问题, 是 NP 难问题。

由于 f(S) 函数具有子模特性, 我们可以采用贪心算法近似求解: 1. 以空集合为起点, 即初始 S=∅; 2. 经过 k 次迭代, 每次选择最大化 f(S+v)-f(S) 的节点 v。贪心算法可以实现至少 (1-1/e) 的近似效果。

信息传播元素: 发送者, 也称作信息源或“种子节点”, 指在信息传递开始时拥有信息的那一小部分用户集合。接收者, 指作为潜在传播目标的广大用户集合。接收者集合的规模要远大于前者, 且不同发送者的目标集合存在重叠。媒介, 指传播过程发生的平台。

1.4 假设有三个城市, 编号分别为 1、2、3。现在有一个商人在三个城市之间来回穿梭, 已知三个城市作为起点的概率分别为 (0.2, 0.4, 0.4)。同时, 这个商人在城市之间旅行或同城停留的跳转概率如下表所示:

城市编号	->1	->2	->3
1	0.5	0.2	0.3
2	0.3	0.5	0.2
3	0.2	0.3	0.5

初始状态概率向量
状态转移概率矩阵

同时, 还知道三座城市各自晴天/雨天的概率如下表所示:

城市编号	1	2	3
晴天概率	0.5	0.4	0.7
雨天概率	0.5	0.6	0.3

观测状态概率矩阵

在某一次旅行中, 商人连续三天观测到的天气状态是 (晴天、雨天、晴天), 请问, 这三天内该商人最有可能的旅行轨迹是什么? 请给出计算过程。

三个城市

记 S 为所有可能状态集合 S = {s₁, s₂, s₃} = {1, 2, 3}。
I = (s₁, s₂, s₃) 是长度为 3 的状态序列, 其对应的观测序列为 W = (w₁, w₂, w₃) = (晴, 雨, 晴)。
A, B, π 分别为隐含状态转移概率矩阵, 观测状态概率矩阵和初始状态概率向量。

(1) 初始化, 在 t = 1 时: 对每个状态 s_i, 计算状态为 s_i, 观测 w₁ 为晴天的概率, 记此概率为 a₁(s_i), 则

代入实际数据得,
a₁(1) = 0.2 × 0.5 = 0.1
a₁(2) = 0.4 × 0.4 = 0.16
a₁(3) = 0.4 × 0.7 = 0.28

x₁ = (4, 1), x₂ = (2, 3), x₃ = (5, 4), x₄ = (1, 0)

$$X = \begin{pmatrix} 4 & -3 & 1 & -2 \\ 2 & -3 & 3 & -2 \\ 5 & -3 & 4 & -2 \\ 1 & -3 & 0 & -2 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ -1 & 1 \\ 2 & 2 \\ -2 & -2 \end{pmatrix}$$

$$R = \begin{bmatrix} 0 & 0 & 0 \\ 0.5 & 0 & 1 \\ 0.5 & 1 & 0 \end{bmatrix}$$

$$A = dR + \left[\frac{(1-d)}{N} \right] e e^T = 0.85 \times \begin{bmatrix} 0 & 0 & 0 \\ 0.5 & 0 & 1 \\ 0.5 & 1 & 0 \end{bmatrix} + 0.05 \times \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 0.05 & 0.05 & 0.05 \\ 0.475 & 0.05 & 0.9 \\ 0.475 & 0.05 & 0.05 \end{bmatrix}$$

$$d_2(s_i) = \max_j \{0.1 \times 0.5, 0.16 \times 0.3, 0.28 \times 0.2\} \times 0.5 = 0.028$$

$$v_2(1) = \operatorname{argmax}_j \{0.1 \times 0.5, 0.16 \times 0.3, 0.28 \times 0.2\} = 3$$

$$d_2(s_i) = \max_j \{0.1 \times 0.2, 0.16 \times 0.5, 0.28 \times 0.3\} \times 0.6 = 0.0504$$

$$v_2(2) = \operatorname{argmax}_j \{0.1 \times 0.2, 0.16 \times 0.5, 0.28 \times 0.3\} = 3$$

$$d_2(s_i) = \max_j \{0.1 \times 0.3, 0.16 \times 0.3, 0.28 \times 0.5\} \times 0.3 = 0.042$$

$$v_2(3) = \operatorname{argmax}_j \{0.1 \times 0.3, 0.16 \times 0.3, 0.28 \times 0.5\} = 3$$

特征值

$$C = X^T X = \begin{pmatrix} 1 & -1 & 2 & -2 \\ -1 & 1 & 2 & -2 \\ 2 & 2 & 2 & -2 \\ -2 & -2 & 2 & -2 \end{pmatrix} = \begin{pmatrix} 10 & 6 \\ 6 & 10 \end{pmatrix}$$

$$\begin{vmatrix} 10-\lambda & 6 \\ 6 & 10-\lambda \end{vmatrix} = 0 \Rightarrow \lambda = 4, 16$$

$$\text{特征向量: } \begin{pmatrix} 10-16 & 6 \\ 6 & 10-16 \end{pmatrix} w_1 = 0 \Rightarrow w_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = w^* = (w_1) \quad X' = X \times W^*$$

1.1 给定以下词项的 idf 值, 以及在三篇文档中的 tf, 已知总文档数为 811,400, 请完成如下计算任务:

	df	tf@Doc1	tf@Doc2	tf@Doc3
Car	18,871	34	8	32
Auto	3,597	3	24	0
Insurance	19,167	0	51	6
Best	40,014	18	0	13

- 计算所有词项的 tf-idf 值。
- 试采用欧式归一化方法 (即向量各元素平方和为 1), 得到处理后的各文档向量化表示, 其中每个向量为 4 维, 每一维对应 1 个词项。
- 基于 2) 中得到的向量化表示, 对于查询“car insurance”, 计算 3 篇文档的得分并进行排序, 其中, 查询中出现的词项权重为 1, 否则为 0。

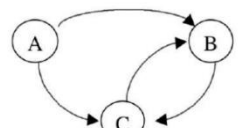
	tf-idf@doc1	TF-idf@doc2	tf-idf@doc3
Car	4.135	3.109	4.092
Auto	3.476	5.601	0
Insurance	0	4.404	2.892
Best	2.947	0	2.763

$$\text{doc1} = \frac{\vec{w}_1}{|\vec{w}_1|} = \left(\frac{4.135}{\sqrt{37.86561}}, \frac{3.476}{\sqrt{37.86561}}, 0, \frac{2.947}{\sqrt{37.86561}} \right) = (0.672, 0.565, 0, 0.478)$$
$$\text{doc2} = (0.400, 0.720, 0.567, 0)$$
$$\text{doc3} = (0.715, 0, 0.505, 0.483)$$

注意归一化方法, 欧式归一

1.2 考虑右图的网络结构图

- 当 Restart 部分的随机跳转概率为 0.15 时, 写出 PageRank 的 (随机) 转移概率矩阵。
- 计算各个节点所对应的 PageRank 值、Hub 值和 Authority 值。



(1) 归一化计算跳转矩阵 R

$$R = \begin{bmatrix} 0 & 0 & 0 \\ 0.5 & 0 & 1 \\ 0.5 & 1 & 0 \end{bmatrix}$$

$$A = dR + \left[\frac{(1-d)}{N} \right] e e^T = 0.85 \times \begin{bmatrix} 0 & 0 & 0 \\ 0.5 & 0 & 1 \\ 0.5 & 1 & 0 \end{bmatrix} + 0.05 \times \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 0.05 & 0.05 & 0.05 \\ 0.475 & 0.05 & 0.9 \\ 0.475 & 0.05 & 0.05 \end{bmatrix}$$

(2) 计算转移矩阵 A

$$A = dR + \left[\frac{(1-d)}{N} \right] e e^T = 0.85 \times \begin{bmatrix} 0 & 0 & 0 \\ 0.5 & 0 & 1 \\ 0.5 & 1 & 0 \end{bmatrix} + 0.05 \times \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$A = \begin{bmatrix} 0.05 & 0.05 & 0.05 \\ 0.475 & 0.05 & 0.9 \\ 0.475 & 0.05 & 0.05 \end{bmatrix}$$

(4) 迭代计算 Pagerank 值

令每个节点的初始 PageRank 值为 1/3

$$\text{由 } P_{n+1} = A P_n \text{ 迭代计算可 } p = \begin{pmatrix} 0.05 \\ 0.475 \\ 0.475 \end{pmatrix}$$

(1) 写出跳转矩阵

$$M = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

(2) 迭代计算并每步进行归一化

$$a_{k+1} = M^T h_k, \quad h_{k+1} = M a_{k+1}$$

$$h = \begin{bmatrix} \sqrt{6}/3 \\ \sqrt{6}/6 \\ \sqrt{6}/6 \end{bmatrix}, \text{authority} = \begin{bmatrix} 0 \\ \sqrt{2}/2 \\ \sqrt{2}/2 \end{bmatrix}$$

归一化