

- 一、问题。① 4个组件指代的部分分别是什么
② WordNet semantic tree 如何构建
③ 3.2中先 S_i 与 β 比较还是先正则化

二、方法。

1° 根据相似度建立一个候选集 (string 上的相似度)

C_i in $O_h^{(m)}$, O_j in $O_e^{(m)}$

① 计算每一个 C_i 与 O_j 的相似度 S_{ij} ($m \times n$)

② * S_i 指 C_i 与 O_e 的相似度, 且若 $S_{ij} > \alpha$, 则

$$S_i = \sum_{j=0}^n S_{ij}$$

③ 若 $S_i > \beta$, 则 C_i 插入候选集 C .

相似度的求法: edit-distance 和 word-net based 2种方法结合.

① edit-distance based

$|fopk|$ 指词 w_i 转换为 w_j 的操作次数

$|w_i|$ 词长

$$e_d(w_i, w_j) = \frac{|fopk|}{\max(|w_i|, |w_j|)}$$

$$S_{ij} = \frac{1}{1 + e_d(i, j)}$$

② word-net based (word-net semantic tree)

v_i 表示 w_i 和 w_j 的第一个共同祖先

v_i, v_j 指 Word i, j 在 word-net tree 中的节点

$count(v_i)$ 指以 v_i 为根节点的子树中节点个数

total 整棵树的总节点数

$$P(v_i) = \frac{count(v_i)}{total}$$

$$S_{ij} = \frac{2 \times \log P(v_i)}{\log P(v_i) + \log P(v_j)}$$

若①或②方法中有 $S_{ij} = 1$, 则 $S_{ij} = 1$, 否则 $S_{ij} = \frac{①+②}{2}$

求 C_i 与 O_e 的 similarity $S_i = \sum_{j=0}^n S_{ij}$; 最后若 $S_i > \beta$, 则 C_i 放入候选集 C

若 $S_{ij} > \alpha$, 则

将 S_i 正则化, 即除以 $\max(S_i)$

2° 基于候选集 C , 根据相关性构造子本体,

(语义上的相似度) 指 concept 对 O_e

相关性受两个因素的影响, 前文的 S_i 和 influence; influence 指 C_i 的邻近结点对 O_e 的相似度和

influence 会随着 distance 而减小, 符合高斯函数的特性. (从图知 x 离 $(0,0)$ 远, 则 $f(x)$ 小)

① $|C_i - C_j|$ 指 C_i 与 C_j 的距离, 即最短路径

C_j 是 C_i 的一个邻近点: 若 $|C_i - C_j| \leq 2$ (在本体图中)

$$\varphi_j(C_i) = S_j \times e^{-|C_i - C_j|^2} \quad (i \neq j)$$

(C_i, C_j 均属于候选集 C)

$$\therefore \varphi(C_i) = \sum \varphi_j(C_i)$$

② C_i 对 O_e 的 relevance 为:

$$r_i = S_i \times \varphi(C_i)$$

③ 若 $r_i < \gamma$, 则将 C_i 从 C 中移除

此时得到子本体所需的 concepts

④ 将 C 中剩下的 concepts 构成有向连接图.

g_i , 且 g_i 是 O_h 的极大连通子图 (sub-graph)

$D(O_e)$ 表示 O_e 的 degree

$$D(O_e) = \frac{\text{边}}{\text{点}}$$

$D(g_i)$ 表示每个 g_i 的度 = 出 + 入度

$|D(O_e) - D(g_i)| > \tau$, 则去掉 g_i 顶点

\therefore degree 差别大, 图的结构相差也大

即不相似 (或)

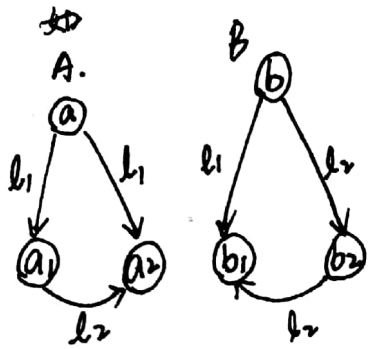
剩下的图 g_i 为子本体.

求 concepts 之间的相似度

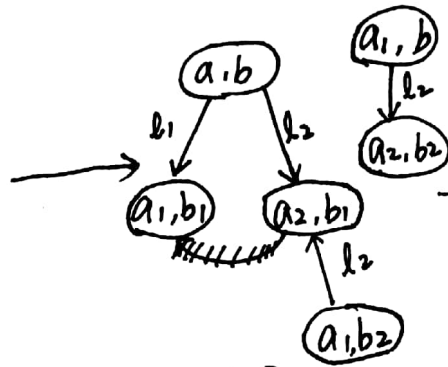
求 concept 与 O_e 的相似度

3. 匹配 O_e 和子本体 \Rightarrow 用 similarity flooding 方法

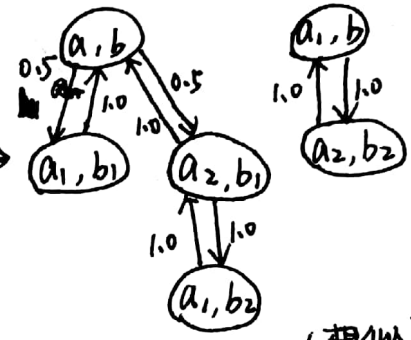
① 将两个本体构建一个有向图, 即图 2.



原本体图 1



连通图 2



(相似度传播图)
induced propagation 图 3

l_1, l_2 指 relation, 若本体 A 和本体 B 的点间拥有同样的 relation, 则变成图 2
 $\langle v_i, v_j \rangle$

② 对于图 2 的每一条边都有权重 w_{ij} , 且 $w_{ij} = \frac{1}{v_i \text{ 出度}}$

③ s_i^0 表示顶点 v_i 的初始相似度, 即第 0 次迭代的相似度.

~~$\sum_j s_j^n \times w_{ij}$~~ 指 v_i 的直接邻点的第 n 次迭代相似度 \times 权重 (从 v_i 离开)
 $\sum_j s_j^n \times w_{ji}$ 指 所有 (进入 v_i)

$$\therefore s_i^{n+1} = s_i^n + \sum_j s_j^n \times w_{ij} + \sum_j s_j^n \times w_{ji}$$

直到 $s_i^{n+1} - s_i^n \leq \alpha$, 则输出 $M(O_e, O_n)$

\therefore 有可能 $s_i^{n+1} - s_i^n > \alpha$, 则去掉