其中, $\alpha$ 是要学习的特征权重, $f(y_i,x_i)$ 是一个概率向量,由各类相似度指标构成(比如文献[26]采用了两个实体在百科中链入(出)实体集的相似度、标签相似度和编写者的相似度等)。

第二类为 PCG 中边的特征函数 $g(y_i, G(y_i))$ , 如公式(6-7)所示,用来表示节点之间的相关性。

$$g(y_i, G(y_i)) = \frac{1}{Z_{\beta}} \exp\{\sum_{y_j \in G(y_i)} \beta g(y_i, y_j)\}$$
 (6-7)

其中, $\beta$ 是要学习的特征权重, $G(y_i)$ 是与 $y_i$ 在 PCG 中的邻居节点集合,函数 $g(y_i,y_j)$ 是一个指示函数,表示在 PCG 中是否存在节点 $x_i$ 到节点 $x_j$ 的边。如图 6-11 所示, $g(y_1,y_3)=1$ 表示 PCG 中的节点 $x_1=(a_2b_3)$ 与节点 $x_3=(a_1b_1)$ 之间存在边,在图中表现为知识图谱 $K_1$ 中实体 $a_1$ 指向了实体 $a_2$ ,在知识图谱 $K_2$ 中实体 $b_1$ 指向了实体 $b_3$ 。

第三类为约束特征函数 $h(y_i, H(y_i))$ ,如公式(6-8)所示,用来表示 PCG 中节点之间的约束。这里的约束指的是一一对应约束,即一个知识图 谱中的实体至多只能和另一个知识图谱中的一个实体等价。

$$h(y_i, H(y_i)) = \frac{1}{Z_{\gamma}} \exp\{\sum_{y_i \in H(y_i)} \gamma h(y_i, y_j)\}$$
(6-8)

其中, $\gamma$ 是要学习的特征权重, $H(y_i)$ 是与 $y_i$ 存在标记冲突(即违反一一对应约束)的隐变量集合, $y_j \in H(y_i)$ 表示两个节点 $x_i$ 和 $x_j$ 中存在一个相同实体, $h(y_i,y_j)$ 是约束函数。当 $y_i = 1$ 并且 $y_j = 1$ 时,也就是说两个节点 $x_i$ 和 $x_j$ 中的实体对都是等价的,因此违反了一一对应约束,此时 $h(y_i,y_j) = 0$ 。在其他情况下, $h(y_i,y_i) = 1$ 。

最后,因子图模型的目标是最大化联合概率P(Y),如公式(6-9)所示,它是三个特征函数的乘积。文献[26]根据给定一组已知匹配的实体对,通过最大化P(Y)似然函数来估计参数 $\alpha$ 、 $\beta$ 和Y的值。根据习得的联合概率模型,实体对的集体对齐问题就转变成根据已对齐实体对估计其他实体对是否能够对齐的条件概率估计问题。如图 6-11 所示的问题就是条件概率  $P(y_1,y_4,y_5|y_2=0,y_3=1)$ 的估计问题。

$$P(Y) = \prod_{i} f(y_{i}, x_{i})g(y_{i}, G(y_{i}))h(y_{i}, H(y_{i}))$$
(6-9)