

关于论文《Discriminative metric learning for multi-view graph partitioning》的学习报告

在论文中,作者们提出了一种学习方法能使初始的多视点结构通过学习最终生成一个稳定的结构。相对于原始的结构,最终生成的稳定结构具有更强的可区分性,从而使得图划分算法和聚类算法能更有效地应用在多视点结构中。在这种学习方法中,作者们假设多视点结构是一个动态的系统结构,在这个结构中,同一视点中存在的节点之间存在相互影响以及不同视点中的节点也会存在相互影响,在更新某一对节点之间的 **connection** 时,与之相关的以上两种作用都应该考虑在内。以下是对算法细节的讨论:

- 1) 关于算法中表征多视点结构中节点之间的连接强度的变量——jaccard similarity,对于无权图,在同一视点中不同节点的连接强度定义为:在该视点中两节点的公共邻居的个数除以两个节点的邻居的总和;同一节点在不同视点中的连接强度定义为:该节点在两视点的公共邻居之和除以该节点在两个视点中的邻居的总和。对于有权图的定义与无权图相似,只是在无权图中统计的是邻居的个数,而在有权图中统计的是邻居与节点的权重和。(无权图相当于把有权图的 **connection** 权重全部置为 1 的特殊情况)。
- 2) 对 **intra-view connection** 的分析,这一部分的影响来源于与节点直接相连的节点的贡献,可以分为三点:第一点是同一节点在不同视点间的作用程度,正如在社交平台中一个用户在某一个平台的行为也会影响该用户在另一个平台的行为;第二点是来自两个节点的公共邻居的影响,假如在同一个社交平台中,两个用户有共同的好友,那么理论上分析他们之间的 **connection** 应该会更强烈;第三点是来自于两个节点的非公共邻居的邻居的影响,这一点与前面两点不一样的是,这个影响并不一定是对 **connection** 产生叠加性影响的,也有可能是产生减弱性影响,假如 p 是 u 的邻居但不是 v 的邻居,但是 p 与 v 的相似性很大,那么这就证明 u 的好友很大程度上都是 v 这种类型的,这样就导致 u 和 v 之间的 **connection** 应该更强,这里定义了一个参数 λ 表征了相似程度的阈值。关于参数 k_{us} 和 c_{us} 的讨论,这两个参数是表征同一视点和跨视点的连接强度的,假如连接强度大,那么这一连接就不容易受其他节点的影响,我有一个猜想是这两个参数的作用是导致算法中的多视点网络结构最终能收敛的原因之一,这是一种类似于负反馈的机制,假如两点之间的连接强度大同时又受其他节点的影响大时,那么这个正反馈最终会导致我们的结构崩溃。
- 3) 对 **inter-view connection** 的分析,这一部分的影响可以分为两点:第一点是两个视点中与节点 u 相关的结构的相似性,假如在两个视点中与节点 u 相关的结构很相似,表明在这两个社交平台中该用户的行为表现都差不多,那么对于该用户来说两个社交平台的连通强度就很大;第二点是来自公共邻居的影响,假如一个用户的朋友在多个社交平台都表现活跃并且这些朋友在各个平台之间的连接强度很大,那么该用户在多个平台之间连接强度大的可能性就很高。
- 4) 关于聚类的准确性的度量,里面用到了一个很重要的函数——香农熵,香农熵是集合信息的度量方式,通过香农熵,我们能计算出以某种方式划分数据(聚类)得到的信息增益,通过信息增益的比较从而能表征出聚类的准确程度。
- 5) 一些感想:从实验结果可以看出这个算法的效果很好,但是我个人认为也很难想得这么周全,如果仅仅是通过类比社交网络中的用户之间的关系从中发现出相互作用的规律想必需要很强的逻辑推理以及类比能力吧。