

HXH 2022.10.28

# 1. Homogeneous Network Embedding 利用 word2vec 应用在同质图

给定一个图  $G = (V, E)$

条件独立性假设

目标: 
$$\arg \max_{\theta} \prod_{v \in V} \prod_{c \in N(v)} P(c|v; \theta) \quad (1)$$

计算时, 把相乘之项转换成对数函数:

$$\arg \max_{\theta} \sum_{v \in V} \sum_{c \in N(v)} \log P(c|v; \theta)$$

$N(v)$  表示节点  $v$  的邻域

$P(c|v; \theta)$  表示给定  $v$  时, 有  $c$  的条件概率

# 2. Heterogeneous Network Embedding (mepach2vec) 处理异质图

## ① Heterogeneous Skip-gram 因为考虑了节点类型

给定:  $G = (V, E, T)$  with  $|T_V| > 1$

$V \rightarrow T_V$   
 $E \rightarrow T_E$

$$|T_V| + |T_E| > 2$$

最大化存在异质上下文  $N_{t,v}$  的概率:

$$\arg \max_{\theta} \sum_{v \in V} \sum_{t \in T_V} \sum_{c_t \in N_{t,v}} \log P(c_t|v; \theta) \quad (2)$$

类型

属于该类型的邻居 (所以说是 Heterogeneous)

$N_{t,v}$ :  $v$  的邻居, 节点类型为第  $t$  种

$$P(c_t|v; \theta) = \frac{e^{X_{c_t} \cdot X_v}}{\sum_{u \in V} e^{X_u \cdot X_v}}$$

注意: 这里未考虑节点类型 ( $t \in T_V$ )  
softmax

$X_v$ : 矩阵  $X$  的第  $v$  行, 表示节点  $v$  的向量表示

注意: 不同类型的节点将会嵌入到同一维度空间

## ② 负采样

给定一个负采样样本量  $M$

(2) 式变成: 
$$\log \sigma(X_{c_t} \cdot X_v) + \sum_{m=1}^M E_{u^m \sim p(u)} [\log \sigma(-X_{u^m} \cdot X_v)]$$

其中  $\sigma(x) = \frac{1}{1+e^{-x}}$

$p(u)$ : 预先定义分布

详细理解见负采样笔记

3. 基于无路径的随机游走 能够确保不同类型的节点之间的关系, 嵌入至 skip-gram 中. 建立每个节点的异构邻域

给定异质网络:  $G = (V, E, T)$

无路径框架:  $\mathcal{D}: \underline{v_1} \xrightarrow{R_1} v_2 \xrightarrow{R_2} \dots v_t \xrightarrow{R_t} v_{t+1} \xrightarrow{R_{t+1}} \dots \underline{v_l}$   $v_1$  与  $v_l$  同一种类型

第  $t$  步的转移概率被定义为:

$$p(v^{t+1} | v^t, \mathcal{D}) = \begin{cases} \frac{1}{|N_{t+1}(v^t)|} & (v^{t+1}, v^t) \in E, \phi(v^{t+1}) = t+1 \\ 0 & (v^{t+1}, v^t) \in E, \phi(v^{t+1}) \neq t+1 \\ 0 & (v^{t+1}, v^t) \notin E \end{cases}$$

邻居节点不是  $t+1$  类型

其中  $v^t \in V_t$ .

$N_{t+1}(v^t)$  表示  $v^t$  节点的邻域中属于  $V_{t+1}$  类型的节点 (根据定义的无路径)

此外: 无路径第一个节点类型与最后一个节点类型一致

如  $p(v^{t+1} | v^t) = p(v^{t+1} | v_1)$ , if  $t=l$

4. metapath2vec++ 同时实现异构网络中结构和语义关联的建模

之前: 给定:  $G = (V, E, T)$ , with  $|T_V| > 1$   
最大化存在异质上下文  $N_t(v)$  的概率:

$$\arg \max_{\theta} \sum_{v \in V} \sum_{t \in T_V} \sum_{c_t \in N_t(v)} \log p(c_t | v; \theta) \quad (2)$$

$N_t(v)$ :  $v$  的邻居, 节点类型为第  $t$  种

$$p(c_t | v; \theta) = \frac{e^{x_{c_t} \cdot x_v}}{\sum_{u \in V} e^{x_u \cdot x_v}} \quad \text{softmax}$$

在 softmax 中未考虑每个节点的类型.

$x_v$ : 矩阵  $X$  的第  $v$  行, 表示节点  $v$  的向量表示

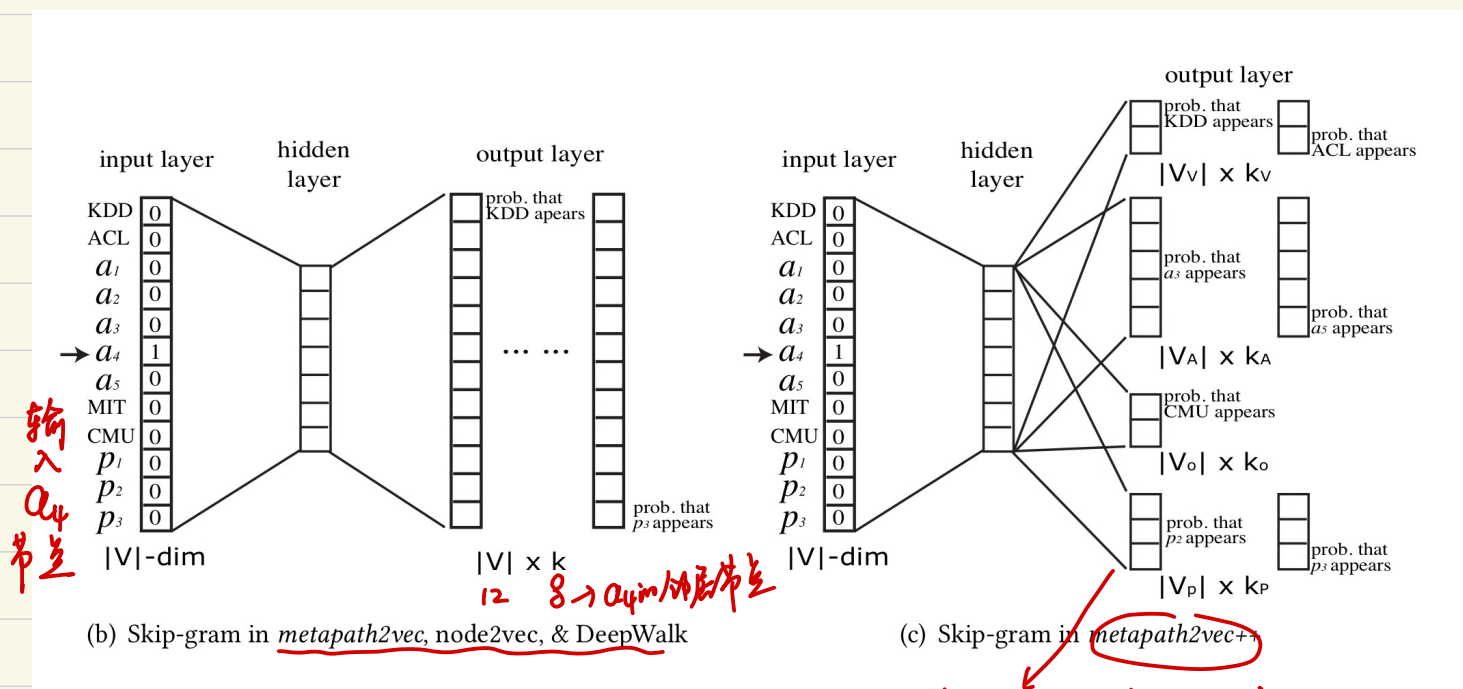
改进: 根据具体的节点类型  $t$  来调整  $p(c_t|v; \theta)$  更新 softmax

$$p(c_t|v; \theta) = \frac{e^{x_{c_t} \cdot x_v}}{\sum_{u_t \in V_t} e^{x_{u_t} \cdot x_v}}$$

也就是说: 在归一化过程中, 会根据固定类型的节点来调整

如此操作, `metapath2vec++` 实际上是在 skip-gram 模型的输出层

考虑了节点  $v$  的每个邻居节点的节点类型



负采样: 没写.

给定一个负采样样本量  $M$

(2) 式更新:  $\log \sigma(x_{c_t} \cdot x_v) + \sum_{m=1}^M E_{u_t^m \sim p(u_t)} [\log \sigma(-x_{u_t^m} \cdot x_v)]$

其中  $\sigma(x) = \frac{1}{1+e^{-x}}$

$p(u_t)$ : 预先定义分布

每种类型对应一个分布

$$O(x) = \log \sigma(x_{c_t} \cdot x_v) + \sum_{m=1}^M E_{u_t^m \sim p_t(u_t)} [\log \sigma(-x_{u_t^m} \cdot x_v)]$$

总结: 比较目标函数

$$\text{metapath2vec} : \log \sigma(X_{c_t} \cdot X_v) + \sum_{m=1}^M E_{u^m \sim p(u)} [\log \sigma(-X_{u^m} \cdot X_v)]$$

$$\text{metapath2vec++} : \log \sigma(X_{c_t} \cdot X_v) + \sum_{m=1}^M E_{u_t^m \sim p_t(u_t)} [\log \sigma(-X_{u_t^m} \cdot X_v)]$$

序列

metapath2vec 异质体现在采样策略上.

{ metapath2vec++ 的异质不仅体现在序列采样策略上, 还体现在负采样.  
(目标函数)