

HXH 2022.10.27

Node2Vec

相比于 DeepWalk, Node2Vec 可以决定去探索节点中的哪一部分.. 也就是有偏的随机游走 \rightarrow 参数 p 和 q

Node2Vec 随机游走机: 完全随机

1. 基本定义

$G = (V, E)$: 给定的图

$f: V \rightarrow \mathbb{R}^d$: 节点到特征表示的映射函数. 也是 Node2Vec 的学习目标. 用于下游的预测任务

d : 参数, 具体化嵌入的维度
所以 f 是一个 $|V| \times d$ 大小的矩阵

对于每一个源节点 $u \in V$, 定义 $N_S(u) \subseteq V$ 为通过一种邻居采样策略 S 采样得到的 u 的邻居节点集合

2. 目标函数

给定节点 u 的向量表示 (通过 f) 的条件下, 最大化其邻域出现的概率 \rightarrow 保证相似性

也是 word2vec 中 skip-gram 模型的思想

$$\max_f \sum_{u \in V} \log \Pr(N_S(u) | f(u))$$

$\rightarrow u$ 邻居节点
 \rightarrow 节点 u embedding
对于每一节点都是如此

作条件独立性假设, 同图节点互不影响 \Rightarrow
假设 ①

$$\Pr(N_S(u) | f(u)) = \prod_{n_i \in N_S(u)} \Pr(n_i | f(u))$$

两个假设都是为了方便处理

→ 节点和其余节点

然后,再假设两个节点之间的相互影响的程度一样(对称性). 假设②

可用 softmax:

$$Pr(n_i | f(u)) = \frac{\exp(f(n_i) \cdot f(u))}{\sum_{u \in V} \exp(f(u) \cdot f(u))}$$

→ 用点积方便计算

根据以上假设, 损失函数可转换为

$$\max_f \sum_{u \in V} \log \prod_{n_i \in N_S(u)} Pr(n_i | f(u))$$

$$= \max_f \sum_{u \in V} \log \prod_{n_i \in N_S(u)} \frac{\exp(f(n_i) \cdot f(u))}{\sum_{u \in V} \exp(f(u) \cdot f(u))}$$

$$= \max_f \sum_{u \in V} \left[\sum_{n_i \in N_S(u)} f(n_i) \cdot f(u) - \log \sum_{u \in V} \exp(f(u) \cdot f(u)) \right]$$

$$= \max_f \sum_{u \in V} \left[-\log Z_u + \sum_{n_i \in N_S(u)} f(n_i) \cdot f(u) \right]$$

其中, $Z_u = \sum_{u \in V} \exp(f(u) \cdot f(u))$ → 称为每个节点的分区函数

→ 为什么要这样处理呢

↓
方便计算了!!!

3. 采样序列策略 → 其实反映了操作者对哪部分信息更重视

网络并不像句子那样是线性的. 因此要生成随机游走序列

① 两种经典(极端)的搜索策略

{ BFS
DFS

如 $N_S(u)$ 不局限于直接
邻居, 还能够广泛地捕获
不同的结构

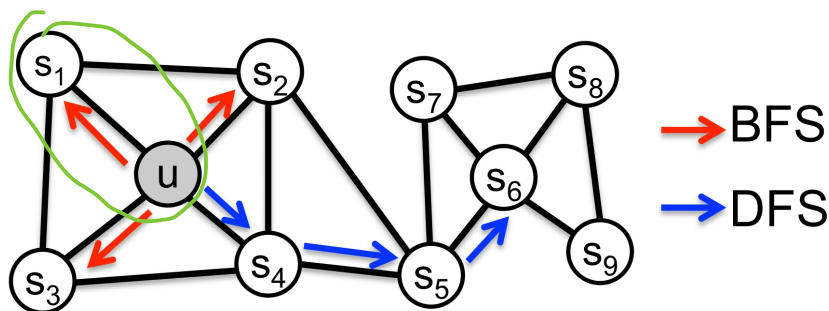


Figure 1: BFS and DFS search strategies from node u ($k = 3$). 数量.

BFS (广度): 搜索的是同质社群 (homophily) (绿色)
 DFS (深度): 搜索的是节点的角色 (中樞、桥接、边缘) (如 u 与 s_6 存在联系)
 但是, 我们需要同时考虑到这两种行为
 所以, Node2Vec 通过设置参数将 BFS 和 DFS 结合

② Node2Vec

2.1 Random Walks

$$P(C_i = x | C_{i-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

π_{vx} → 节点 v 跳到节点 x 的概率
 Z → 归一化常数
 if $(v, x) \in E$ → 表示有连接

2.2 Search bias α (通过 p 和 q 这两个参数定义一个二阶随机游走)

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases}$$

$d_{tx} = \{0, 1, 2\}$ 表示 t 和 x 的最短路径距离

$$\pi_{vx} = \alpha_{pq}(t, x) \cdot w_{vx}$$

w_{vx} → v 与 x 之间的权重

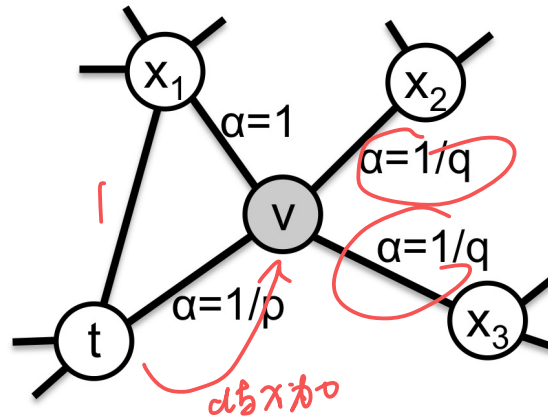


Figure 2: Illustration of the random walk procedure in node2vec. The walk just transitioned from t to v and is now evaluating its next step out of node v . Edge labels indicate search biases α .

刚刚由 $t \rightarrow v$, 现在在 v 节点-考虑下一个

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \text{ } \rightarrow \text{所以 } d_{tx} = 0 \text{ } t \text{ 到 } v \text{ 返回 (} t \text{ 与 } x \text{ 相等)} \\ 1 & \text{if } d_{tx} = 1 \text{ } t \text{ 到 } x_1 \text{ (} t \text{ 与 } x \text{ 相连)} \\ \frac{1}{q} & \text{if } d_{tx} = 2 \text{ } t \text{ 到 } x_2, x_3 \text{ 远行 (} t \text{ 与 } x \text{ 不相连)} \end{cases}$$

$$\tilde{r}_{vx} = \alpha_{pq}(t, x) \cdot W_{vx}$$

参数 p 和 q 的意义:

① 返回概率 p

$\begin{cases} p > \max(q, 1), \text{ 采样尽量不返回, 对应上图, 即下一节点不太可能是上一个节点 } t \\ p < \max(q, 1), \text{ 采样倾向于返回上一节点, 这样会一直在 } x \text{ 节点周围转.} \end{cases}$

② 出入参数 q

$$\begin{cases} q > 1, \text{ 游走倾向于在起始点周围转 (BFS)} \\ q < 1, \text{ 远走 (DFS)} \end{cases}$$

当 $p=1, q=1$, NodeVec 中的游走等同于 DeepWalk 中的随机游走。
也就是说, DeepWalk 是 NodeVec 中的一个特例

2.3 时空复杂度:

时间: $O(\frac{1}{k} + \frac{1}{l-k})$, k 是邻域节点的个数, 整个序列长为 l
所以一个节点已被采样的概率为 $\frac{1}{k}$, 在另一边为 $\frac{1}{l-k}$.

空间: 存储每个节点的邻居的复杂度为 $O(|E|)$

对于 $= \Pi$, 要存储邻居的邻居, 故为 $O(a^2|V|)$

a 表示每个节点的平均连接数, (上一步是 a 个, 当前也是)

4. Node2vec 算法

Algorithm 1 The node2vec algorithm.

LearnFeatures (Graph $G = (V, E, W)$, Dimensions d , Walks per node r , Walk length l , Context size k , Return p , In-out q)

① $\pi = \text{PreprocessModifiedWeights}(G, p, q)$
 $G' = (V, E, \pi)$ 生成随机游走采样策略
Initialize $walks$ to Empty
② **for** $iter = 1$ **to** r **do** } 每个节点生成 r 个随机游走序列
 for all nodes $u \in V$ **do**
 $walk = \text{node2vecWalk}(G', u, l)$ 生成 1 个随机游走序列
 Append $walk$ to $walks$
③ $f = \text{StochasticGradientDescent}(k, d, walks)$ skip-gram 训练得到节点
 return f 嵌入表示

node2vecWalk (Graph $G' = (V, E, \pi)$, Start node u , Length l)

Initialize $walk$ to $[u]$

for $walk_iter = 1$ **to** l **do**

$curr = walk[-1]$

$V_{curr} = \text{GetNeighbors}(curr, G')$

$s = \text{AliasSample}(V_{curr}, \pi)$

Append s to $walk$

return $walk$

理解:

通过 node2vecWalk 获得序列 $walks \Rightarrow$ 放入 skip-gram 中训练
(每个节点)

\Rightarrow 不断优化 \Rightarrow 得到节点的嵌入表示 \rightarrow 看代码再理解.