

# 表示学习之node2vec

原创 不懂技术的天才 人工智能不要人工 2020-09-15

近年来，知识表示学习在学习自身知识特征进行自动预测领域取得了重大进步。然而，现在的特征学习方法还不足以完全表示网络中已有连接模式的多样性。node2vec为网络节点学习连续的特征表示，通过最大化保留邻接节点的概率，将节点映射到低维的特征空间进行向量表示。

在node2vec之前已经有DeepWalk将word2vec思想运用到图的向量表示。node2vec的主要创新点就是在寻找节点邻域（即节点的上下文信息）的方法做了改进。定义一个灵活的节点的邻域概念和设计有偏的随机游走算法来高效探索邻域的多样性。

node2vec的本质思想来源于word2vec。在给定一个节点，然后通过图搜索算法获取节点的邻域集合，用邻域模拟词的上下文信息。最后使用skip-gram模型训练得到节点的向量。

node2vec将特征表示学习定义为一个最大似然优化问题，定义了一个目标优化函数：

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

为了简化优化问题，作了两个假设：

1、条件独立，在给定源节点的特征表示后，观察一个邻域节点的可能性独立于观察任何其他邻域节点。

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

2、特征空间对称性，源节点和近邻节点在特征空间中具有对称效应。

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

$N_S(u)$ 表示源节点 $u$ 的邻域(即邻接节点的集合)， $n_i$ 表示源节点邻域中的第 $i$ 个节点， $f$ 是将节点映射特征表示的函数， $V$ 是网络图的节点集合， $v$ 是网络图中除了源节点的其他节点。**注：符号“ $\cdot$ ”表示后面的字符为下标，对应公式中的表示，下同。**

node2vec的特征学习是基于skip-gram模型的。而skip-gram模型一开始用来处理自然语言文本的，文本句子是线性的。在给定一段文本，某个词的上下文信息直接使用滑动窗口来获取。但是网络图是非线性的，需要为每个节点定义上下文信息的概念。这里采用的是随机过程采样。

随机过程采样的实现是运用在BFS和DFS的基础上设计有偏的随机游走的方式。假设源节点为 $c_0 = u$ ，模拟固定长度 $l$ 的随机游走的过程中，第 $i$ 个节点 $c_i$ 的概率可以表示为：

$$P(c_i = x \mid c_{i-1} = v) = \begin{cases} \frac{\pi_{v,x}}{Z} & \text{if } (v, x) \in E \\ 0 & \text{otherwise} \end{cases}$$

$\pi_{vx}$ 为节点 $v$ 与 $x$ 之间的未归一化转移概率， $Z$ 为归一化常数。

对于有偏变量 $\alpha$ 的计算，定义带有两个参数 $p$ 和 $q$ 的二阶随机游走。假目前一个随机游走走过边 $(t,v)$ ，目前在节点 $v$ 出，现在需要考虑下一个节点往哪走，如图所示：

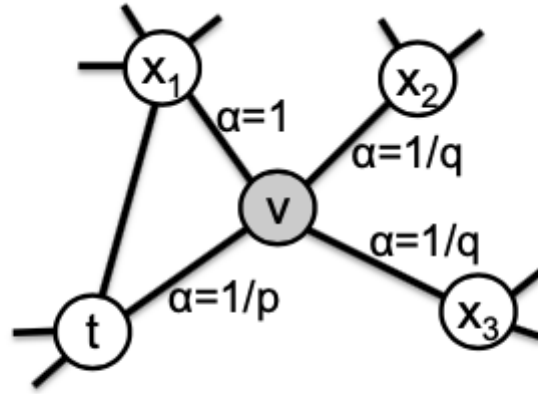


Figure 2: Illustration of the random walk procedure in *node2vec*. The walk just transitioned from  $t$  to  $v$  and is now evaluating its next step out of node  $v$ . Edge labels indicate search biases  $\alpha$ .

因此，需要评估随机游走的转移未归一化概率 $\pi_{vx}$ （即从 $v$ 游走到下一个节点 $x$ 的概率）。定义 $\pi_{vx} = \alpha_{pq}(t,x) \cdot w_{vx}$ ，其中 $\alpha_{pq}(t,x)$ 为边 $(t,x)$ 的偏置项， $w_{vx}$ 为边的权重，若没有权重则 $w_{vx}=1$ 。表达式如图：

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

$d_{tx}$ 表示 $t$ 到 $x$ 的最短距离。且取值范围为 $\{0,1,2\}$ 。从公式可以看出：

1. 参数 $p$ 是控制下一步是否需要重新访问上一个节点。如果 $p$ 值设置的很大（大于 $\max(q,1)$ ），则将基本上不会访问 $t$ 节点。相反，若 $p$ 值设置很小（小于 $\min(q,1)$ ），则很大可能重新访问 $t$ 节点；
2. 参数 $q$ 是控制下一步游走方向是靠近节点 $t$ 的节点和是远离节点 $t$ 的节点。如果 $q$ 设置为1，则下一步将访问 $x_1$ ，如果 $q$ 设置小于1，则下一步将访问 $x_2$ 或 $x_3$ 。也相当于控制随机游走的方式是按广度优先搜索和深度优先搜索的方法进行。

**采用随机游走的优点：随机游走在空间上和时间上的计算复杂度都非常高效；**

最后，*node2vec*算法伪代码如下：

**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  $q$ )  
 $\pi = \text{PreprocessModifiedWeights}(G, p, q)$   
 $G' = (V, E, \pi)$   
Initialize *walks* to Empty  
**for**  $iter = 1$  **to**  $r$  **do**  
    **for all** nodes  $u \in V$  **do**  
         $walk = \text{node2vecWalk}(G', u, l)$   
        Append *walk* to *walks*  
 $f = \text{StochasticGradientDescent}(k, d, walks)$   
**return**  $f$

---

**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*

---

整个过程分三个阶段：边的转移概率的预处理计算、模拟随机游走和使用SGD优化。且每个阶段都可以同步执行。

node2vec源码：

```

1  import numpy as np
2  import networkx as nx
3  import random
4
5
6  class Graph():
7      def __init__(self, nx_G, is_directed, p, q):
8          self.G = nx_G
9          self.is_directed = is_directed
10         self.p = p
11         self.q = q
12
13     def node2vec_walk(self, walk_length, start_node):
14         '''
15         Simulate a random walk starting from start node.
16         '''

```

```
17     G = self.G
18     alias_nodes = self.alias_nodes
19     alias_edges = self.alias_edges
20
21     walk = [start_node]
22
23     while len(walk) < walk_length:
24         cur = walk[-1]
25         cur_nbrs = sorted(G.neighbors(cur))
26         if len(cur_nbrs) > 0:
27             if len(walk) == 1:
28                 walk.append(cur_nbrs[alias_draw(alias_nodes[cur][0], alias_nodes[
29             else:
30                 prev = walk[-2]
31                 next = cur_nbrs[alias_draw(alias_edges[(prev, cur)][0],
32                     alias_edges[(prev, cur)][1])]
33                 walk.append(next)
34         else:
35             break
36
37     return walk
38
39 def simulate_walks(self, num_walks, walk_length):
40     '''
41     Repeatedly simulate random walks from each node.
42     '''
43     G = self.G
44     walks = []
45     nodes = list(G.nodes())
46     print 'Walk iteration:'
47     for walk_iter in range(num_walks):
48         print str(walk_iter+1), '/', str(num_walks)
49         random.shuffle(nodes)
50         for node in nodes:
51             walks.append(self.node2vec_walk(walk_length=walk_length, start_node=
52
53     return walks
54
55 def get_alias_edge(self, src, dst):
56     '''
```

```
57     Get the alias edge setup lists for a given edge.
58     '''
59     G = self.G
60     p = self.p
61     q = self.q
62
63     unnormalized_probs = []
64     for dst_nbr in sorted(G.neighbors(dst)):
65         if dst_nbr == src:
66             unnormalized_probs.append(G[dst][dst_nbr]['weight']/p)
67         elif G.has_edge(dst_nbr, src):
68             unnormalized_probs.append(G[dst][dst_nbr]['weight'])
69         else:
70             unnormalized_probs.append(G[dst][dst_nbr]['weight']/q)
71     norm_const = sum(unnormalized_probs)
72     normalized_probs = [float(u_prob)/norm_const for u_prob in unnormalized_probs]
73
74     return alias_setup(normalized_probs)
75
76 def preprocess_transition_probs(self):
77     '''
78     Preprocessing of transition probabilities for guiding the random walks.
79     '''
80     G = self.G
81     is_directed = self.is_directed
82
83     alias_nodes = {}
84     for node in G.nodes():
85         unnormalized_probs = [G[node][nbr]['weight'] for nbr in sorted(G.neighbors(node))]
86         norm_const = sum(unnormalized_probs)
87         normalized_probs = [float(u_prob)/norm_const for u_prob in unnormalized_probs]
88         alias_nodes[node] = alias_setup(normalized_probs)
89
90     alias_edges = {}
91     triads = {}
92
93     if is_directed:
94         for edge in G.edges():
95             alias_edges[edge] = self.get_alias_edge(edge[0], edge[1])
96     else:
```

```
97     for edge in G.edges():
98         alias_edges[edge] = self.get_alias_edge(edge[0], edge[1])
99         alias_edges[(edge[1], edge[0])] = self.get_alias_edge(edge[1], edge[0])
100
101     self.alias_nodes = alias_nodes
102     self.alias_edges = alias_edges
103
104     return
105
106
107 def alias_setup(probs):
108     '''
109     Compute utility lists for non-uniform sampling from discrete distributions.
110     Refer to https://hips.seas.harvard.edu/blog/2013/03/03/the-alias-method-efficient-sampling-with-many-discrete-outcomes/
111     for details
112     '''
113     K = len(probs)
114     q = np.zeros(K)
115     J = np.zeros(K, dtype=np.int)
116
117     smaller = []
118     larger = []
119     for kk, prob in enumerate(probs):
120         q[kk] = K*prob
121         if q[kk] < 1.0:
122             smaller.append(kk)
123         else:
124             larger.append(kk)
125
126     while len(smaller) > 0 and len(larger) > 0:
127         small = smaller.pop()
128         large = larger.pop()
129
130         J[small] = large
131         q[large] = q[large] + q[small] - 1.0
132         if q[large] < 1.0:
133             smaller.append(large)
134         else:
135             larger.append(large)
136
```

```
137     return J, q
138
139 def alias_draw(J, q):
140     '''
141     Draw sample from a non-uniform discrete distribution using alias sampling
142     '''
143     K = len(J)
144
145     kk = int(np.floor(np.random.rand()*K))
146     if np.random.rand() < q[kk]:
147         return kk
148     else:
149         return J[kk]
```

本文的内容主要是阅读node2vec论文《node2vec: Scalable Feature Learning for Networks》的阅读记录。如有错误欢迎留言指正。

喜欢此内容的人还喜欢

方李莉|谁在成长

中国艺术人类学学会

---

一副流传极广的对联，堪称天下“神联”，中老年人有共鸣

新退休生活