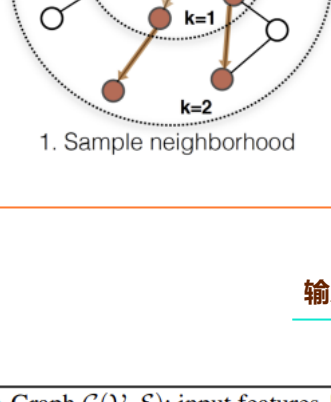


# Graph sage

## 理论

graphsage是可以运用在无监督学习即节点聚类任务的一种方式，因为其sample和aggregate使得它能够具有捕捉图结构的特性。当然，监督任务，即节点分类，边连接预测也是可以做到的。



1. Sample neighborhood



2. Aggregate feature information from neighbors

## 输入参数

**Input** : Graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ ; input features  $\{x_v, \forall v \in \mathcal{V}\}$ ; depth  $K$ ; weight matrices  $W^k, \forall k \in \{1, \dots, K\}$ ; non-linearity  $\sigma$ ; differentiable aggregator functions  $\text{AGGREGATE}_k, \forall k \in \{1, \dots, K\}$ ; neighborhood function  $\mathcal{N}: v \rightarrow 2^{\mathcal{V}}$

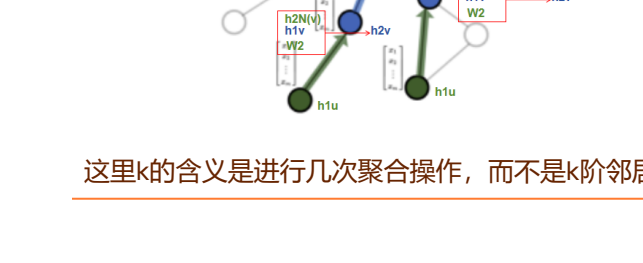
- Graph (V,E):V-节点集合, E-边集合
- input features (xv,  $\forall v \in V$ ):输入的特征
- depth K: 是表示深度计算k阶邻居
- $W_k$  (weight matrices) -指的是对k层在其进行聚合的时候都赋予一个单独的权重矩阵
- neighborhood function  $N: v \rightarrow 2^V$ 表示的是v节点的邻居节点

## 输出参数

**Output**: Vector representations  $z_v$  for all  $v \in \mathcal{V}$

节点v的向量表示

## 实现流程



这里k的含义是进行几次聚合操作，而不是k阶邻居。

## 代码实现

### 1. 数据预处理

#### 1-1. 节点特征矩阵构建、标签列表、节点索引、标签索引

```
from collections import defaultdict
import numpy as np
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.nn.init as init

# 节点特征矩阵构建
def load_data(path):
    with open(path, 'r') as f:
        lines = f.readlines()
        n_nodes = int(lines[0].strip())
        n_labels = int(lines[1].strip())
        features = []
        labels = []
        for i in range(2, len(lines)):
            line = lines[i].strip()
            node_id, feat, label = line.split()
            features.append([float(x) for x in feat.split()])
            labels.append(int(label))
        features = np.array(features)
        labels = np.array(labels)
    return features, labels
```

#### 1-2 构建邻接字典

```
def build_adj_dict(path):
    with open(path, 'r') as f:
        lines = f.readlines()
        n_nodes = int(lines[0].strip())
        n_labels = int(lines[1].strip())
        adj_dict = defaultdict(list)
        for i in range(2, len(lines)):
            line = lines[i].strip()
            node_id, feat, label = line.split()
            adj_dict[node_id].append(int(feat))
        adj_dict = dict(adj_dict)
    return adj_dict
```

### 2. 划分数据集

```
def split_data(features, labels):
    n_nodes = features.shape[0]
    n_labels = labels.shape[0]
    train_size = int(n_nodes * 0.8)
    test_size = n_nodes - train_size
    train_features = features[:train_size]
    train_labels = labels[:train_size]
    test_features = features[train_size:]
    test_labels = labels[train_size:]
    return train_features, train_labels, test_features, test_labels
```

#### 2. 模型架构

##### 2-1 Sage layer

##### 2-1-1 初始化设定

```
class SageLayer(nn.Module):
    def __init__(self, num_layers, input_size, out_size, raw_features, adj_lists, device,
                 gnn_type, agg_func=None):
        super(SageLayer, self).__init__()
        self.num_layers = num_layers
        self.gnn_type = gnn_type
        self.agg_func = agg_func
        self.device = device
        self.raw_features = raw_features
        self.adj_lists = adj_lists
        self.out_size = out_size
        self.in_size = input_size
        self.layers = nn.ModuleList()
        for i in range(self.num_layers):
            self.layers.append(SageLayer(self.in_size, self.out_size, self.device,
                                         self.gnn_type, self.agg_func))
            self.in_size = self.out_size
```

```
def forward(self, self_feats, aggregate_feats, neighbors):
    # Generates embeddings for a batch of nodes.
    nodes = list(nodes)
    if not self.gnn:
        combined = torch.cat([self_feats, aggregate_feats], dim=1)
    else:
        combined = aggregate_feats
    combined = F.relu(self.weight.mul(combined.t()).t())
    return combined
```

##### 2-1-2-1 获取节点的邻居——get\_unique\_neighs\_list

##### 2-2 graphsage

##### 2-2-1 初始化设定

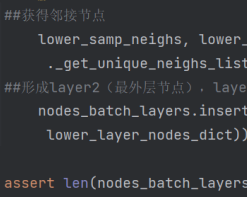
```
class GraphSage(nn.Module):
    """Base class for GraphSage"""
    def __init__(self, num_layers, input_size, out_size, raw_features, adj_lists, device,
                 gnn_type, agg_func=None):
        super(GraphSage, self).__init__()
        self.num_layers = num_layers
        self.gnn_type = gnn_type
        self.agg_func = agg_func
        self.device = device
        self.raw_features = raw_features
        self.adj_lists = adj_lists
        self.out_size = out_size
        self.in_size = input_size
        self.layers = nn.ModuleList()
        for i in range(self.num_layers):
            self.layers.append(SageLayer(self.in_size, self.out_size, self.device,
                                         self.gnn_type, self.agg_func))
            self.in_size = self.out_size
```

##### 2-2-2 前向传播

##### 2-2-2-1 训练批次的节点导入& 聚合层列表的建立

```
def forward(self, nodes_batch):
    # Generates embeddings for a batch of nodes.
    lower_layer_nodes = list(nodes_batch)
    # 聚合层列表
    nodes_batch_layers = [lower_layer_nodes]
```

##### 2-2-2-2 构建采样层列表



1. 先对训练batch的节点 (layer\_center)使用采样邻居函数，获得layer\_center的邻居节点(layer1)随后将其插入nodes\_batch\_layers列表的0号位置，即生成[layer1,layer\_center]

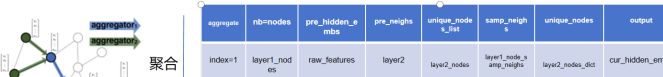
2. 再对layer1节点使用采样邻居函数，获得layer1中节点的邻居节点(layer2) 随后将其插入nodes\_batch\_layers列表的0号位置即生成[layer2,layer1,layer\_center]

3. 总结：通过这样的方式其实是实现了从内到外的采样节点，并且为之后的聚合节点提供了一个聚合层

##### 运用采样邻居函数

```
def get_unique_neighs_list(self, nodes, num_sample=10):
    # Get a list of unique neighbors for a node.
    # nodes: list of nodes
    # num_sample: number of neighbors to sample
    # Returns: list of unique neighbors
    # 1. 获取节点的邻居
    neighbors = []
    for node in nodes:
        neighbors.append(self.get_neighs(node))
    # 2. 去重
    unique_neighs = list(set(neighbors))
    # 3. 采样
    unique_neighs = self.sample(unique_neighs, num_sample)
    return unique_neighs
```

##### 2-2-2-3 聚合邻居节点更新节点信息



聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

聚合

##### 聚合函数

```
def aggregate(self, nodes, pre_hidden_emb, pre_neighs):
    # Aggregate features from neighbors and the node itself.
    # nodes: list of nodes
    # pre_hidden_emb: pre-hidden embeddings
    # pre_neighs: pre-neighbors
    # Returns: aggregated features
    # 1. 聚合
    aggregated_feats = []
    for node in nodes:
        aggregated_feats.append(self.aggregate(node, pre_hidden_emb, pre_neighs))
    return aggregated_feats
```

##### 选取聚合方式进行聚合

```
def aggregate(self, nodes, pre_hidden_emb, pre_neighs):
    # Aggregate features from neighbors and the node itself.
    # nodes: list of nodes
    # pre_hidden_emb: pre-hidden embeddings
    # pre_neighs: pre-neighbors
    # Returns: aggregated features
    # 1. 聚合
    aggregated_feats = []
    for node in nodes:
        aggregated_feats.append(self.aggregate(node, pre_hidden_emb, pre_neighs))
    return aggregated_feats
```

##### sage\_layer

```
class SageLayer(nn.Module):
    def __init__(self, num_layers, input_size, out_size, raw_features, adj_lists, device,
                 gnn_type, agg_func=None):
        super(SageLayer, self).__init__()
        self.num_layers = num_layers
        self.gnn_type = gnn_type
        self.agg_func = agg_func
        self.device = device
        self.raw_features = raw_features
        self.adj_lists = adj_lists
        self.out_size = out_size
        self.in_size = input_size
        self.layers = nn.ModuleList()
        for i in range(self.num_layers):
            self.layers.append(SageLayer(self.in_size, self.out_size, self.device,
                                         self.gnn_type, self.agg_func))
            self.in_size = self.out_size
```

在第二次聚合的时候会使用，主要是为了提取第一次聚合后的节点特征作为其送入sagelayer的初始特征

##### 分类器定义

```
class Classification(nn.Module):
    def __init__(self, emb_size, num_classes):
        super(Classification, self).__init__()
        self.weight = nn.Parameter(torch.FloatTensor(emb_size, num_classes))
        self.layer = nn.Sequential(
            nn.Linear(emb_size, num_classes),
            nn.Softmax()
        )
        self.init_params()

    def init_params(self):
        for param in self.parameters():
            if len(param.size()) == 1:
                nn.init.xavier_uniform_(param)

    def forward(self, emb_size):
        logits = torch.nn.functional.linear(emb_size, self.weight)
        return logits
```

设定一个全连接层，输出部分对应着分类的种类，随后再用一个log\_softmax输出其属于每个分类的概率

##### 2-3-2 损失计算

```
def train(self, loader):
    # Training loop
    # loader: DataLoader
    # Returns: loss
    # 1. 训练
    loss = 0
    for batch in loader:
        # 2. 计算损失
        loss_batch = self.compute_loss(batch)
        loss += loss_batch
    return loss
```

负对数似然损失，通常称为交叉熵损失，它通常用作具有softmax 输出的监督分类任务的损失函数。 logits 是神经网络对每个类别的原始分数（未归一化概率）。

负对数似然损失使用 logits和真实标签 labels, batch 来计算，它衡量了预测的概率与真实类别标签的匹配程度。

最后，将损失在批次中平均，得到每个样本的平均损失，并将这个值赋给变量 loss。