

节点预测与边预测任务实践

引言

在此小节我们将利用在上一小节[6-1-数据完整存于内存的数据集类](#)中构造的 PlanetoidPubMed 数据集类，来实践节点预测与边预测任务。

注：边预测任务实践中的代码来源于[link_pred.py](#)。

节点预测任务实践

之前我们学习过由2层 GATConv 组成的图神经网络，现在我们重定义一个 GAT图神经网络，使其能够通过参数来定义 GATConv 的层数，以及每一层 GATConv 的 out_channels。我们的图神经网络定义如下：

```
1 class GAT(torch.nn.Module):
2     def __init__(self, num_features,
3         hidden_channels_list, num_classes):
4         super(GAT, self).__init__()
5         torch.manual_seed(12345)
6         hns = [num_features] + hidden_channels_list
7         conv_list = []
8         for idx in range(len(hidden_channels_list)):
9             conv_list.append((GATConv(hns[idx],
10                 hns[idx+1]), 'x, edge_index -> x'))
11             conv_list.append(ReLU(inplace=True),)
12
13         self.convseq = Sequential('x, edge_index',
14             conv_list)
15         self.linear = Linear(hidden_channels_list[-1],
16             num_classes)
17
18     def forward(self, x, edge_index):
19         x = self.convseq(x, edge_index)
20         x = F.dropout(x, p=0.5, training=self.training)
21         x = self.linear(x)
22         return x
```

由于我们的神经网络由多个 `GATConv` 顺序相连而构成，因此我们使用了 `torch_geometric.nn.Sequential` 容器，详细内容可见于[官方文档](#)。

我们通过 `hidden_channels_list` 参数来设置每一层 `GATConv` 的 `outchannel`，所以 `hidden_channels_list` 长度即为 `GATConv` 的层数。通过修改 `hidden_channels_list`，我们就可构造出不同的图神经网络。

完整的代码可见于 `codes/node_classification.py`。请小伙伴们自行完成代码中图神经网络类的训练、验证和测试。

边预测任务实践

边预测任务，目标是预测两个节点之间是否存在边。拿到一个图数据集，我们有节点属性 `x`，边端点 `edge_index`。`edge_index` 存储的便是正样本。为了构建边预测任务，我们需要生成一些负样本，即采样一些不存在边的节点对作为负样本边，正负样本数量应平衡。此外要将样本分为训练集、验证集和测试集三个集合。

PyG中为我们提供了现成的采样负样本边的方法，`train_test_split_edges(data, val_ratio=0.05, test_ratio=0.1)`，其

- 第一个参数为 `torch_geometric.data.Data` 对象，
- 第二参数为验证集所占比例，
- 第三个参数为测试集所占比例。

该函数将自动地采样得到负样本，并将正负样本分成训练集、验证集和测试集三个集合。它用 `train_pos_edge_index`、`train_neg_adj_mask`、`val_pos_edge_index`、`val_neg_edge_index`、`test_pos_edge_index` 和 `test_neg_edge_index`，六个属性取代 `edge_index` 属性。

注意 `train_neg_adj_mask` 与其他属性格式不同，其实该属性在后面并没有派上用场，后面我们仍然需要进行一次训练集负样本采样。

下面我们使用 `cora` 数据集作为例子，进行边预测任务说明。

获取数据集并进行分析

首先是获取数据集并进行分析：

```
1 import os.path as osp
2
3 from torch_geometric.utils import negative_sampling
4 from torch_geometric.datasets import Planetoid
5 import torch_geometric.transforms as T
6 from torch_geometric.utils import
  train_test_split_edges
7
8 dataset = Planetoid('dataset', 'Cora',
  transform=T.NormalizeFeatures())
9 data = dataset[0]
10 data.train_mask = data.val_mask = data.test_mask =
  data.y = None # 不再有用
11
12 print(data.edge_index.shape)
13 # torch.Size([2, 10556])
14
15 data = train_test_split_edges(data)
16
17 for key in data.keys:
18     print(key, getattr(data, key).shape)
19
20 # x torch.Size([2708, 1433])
21 # val_pos_edge_index torch.Size([2, 263])
22 # test_pos_edge_index torch.Size([2, 527])
23 # train_pos_edge_index torch.Size([2, 8976])
24 # train_neg_adj_mask torch.Size([2708, 2708])
25 # val_neg_edge_index torch.Size([2, 263])
26 # test_neg_edge_index torch.Size([2, 527])
27 # 263 + 527 + 8976 = 9766 != 10556
28 # 263 + 527 + 8976/2 = 5278 = 10556/2
```

我们观察到训练集、验证集和测试集中正样本边的数量之和不等于原始边的数量。这是因为，现在所用的Cora图是无向图，在统计原始边数量时，每一条边的正向与反向各统计了一次，训练集也包含边的正向与反向，但验证集与测试集都只包含了边的一个方向。

为什么训练集要包含边的正向与反向，而验证集与测试集都只包含了边的一个方向？这是因为，训练集用于训练，训练时一条边的两个端点要互传信息，只考虑一个方向的话，只能由一个端点传信息给另一个端点，而验证集与测试集的边用于衡量检验边预测的准确性，只需考虑一个方向的边即可。

边预测图神经网络的构造

接下来构造神经网络：

```
1 import torch
2 from torch_geometric.nn import GCNConv
3
4 class Net(torch.nn.Module):
5     def __init__(self, in_channels, out_channels):
6         super(Net, self).__init__()
7         self.conv1 = GCNConv(in_channels, 128)
8         self.conv2 = GCNConv(128, out_channels)
9
10    def encode(self, x, edge_index):
11        x = self.conv1(x, edge_index)
12        x = x.relu()
13        return self.conv2(x, edge_index)
14
15    def decode(self, z, pos_edge_index,
16               neg_edge_index):
17        edge_index = torch.cat([pos_edge_index,
18                                neg_edge_index], dim=-1)
19        return (z[edge_index[0]] *
20                z[edge_index[1]]).sum(dim=-1)
21
22    def decode_all(self, z):
23        prob_adj = z @ z.t()
24        return (prob_adj >
25                0).nonzero(as_tuple=False).t()
```

用于做边预测的神经网络主要由两部分组成：其一是编码（encode），它与我们前面介绍的节点表征生成是一样的；其二是解码（decode），它根据边两端节点的表征生成边为真的几率（odds）。`decode_all(self, z)` 用于推理（inference）阶段，我们要对所有的节点对预测存在边的几率。

边预测图神经网络的训练

定义单个epoch的训练过程

```
1 def get_link_labels(pos_edge_index, neg_edge_index):
2     num_links = pos_edge_index.size(1) +
neg_edge_index.size(1)
3     link_labels = torch.zeros(num_links,
dtype=torch.float)
4     link_labels[:pos_edge_index.size(1)] = 1.
5     return link_labels
6
7 def train(data, model, optimizer):
8     model.train()
9
10    neg_edge_index = negative_sampling(
11        edge_index=data.train_pos_edge_index,
12        num_nodes=data.num_nodes,
13
14        num_neg_samples=data.train_pos_edge_index.size(1))
15
16    optimizer.zero_grad()
17    z = model.encode(data.x, data.train_pos_edge_index)
18    link_logits = model.decode(z,
data.train_pos_edge_index, neg_edge_index)
19    link_labels =
get_link_labels(data.train_pos_edge_index,
neg_edge_index).to(data.x.device)
20    loss =
F.binary_cross_entropy_with_logits(link_logits,
link_labels)
21    loss.backward()
22    optimizer.step()
23
24    return loss
```

通常，存在边的节点对的数量往往少于不存在边的节点对的数量。我们在每一个 epoch 的训练过程中，都进行一次训练集负样本采样。采样到的样本数量与训练集正样本相同，但不同 epoch 中采样到的样本是不同的。这样做，我们既能实现类别数量平衡，又能实现增加训练集负样本的多样性。在负样本采样时，我们传递了 `train_pos_edge_index` 为参数，于是 `negative_sampling()` 函数只会在训练集中不存在边的节点对中采样。`get_link_labels()` 函数用于生成完整训练集的标签。

注：在训练阶段，我们应该只见训练集，对验证集与测试集都是不可见的。所以我们没有使用所有的边，而是只用了训练集正样本边。

定义单个epoch验证与测试过程

```
1  @torch.no_grad()
2  def test(data, model):
3      model.eval()
4
5      z = model.encode(data.x, data.train_pos_edge_index)
6
7      results = []
8      for prefix in ['val', 'test']:
9          pos_edge_index =
10             data[f'{prefix}_pos_edge_index']
11             neg_edge_index =
12             data[f'{prefix}_neg_edge_index']
13             link_logits = model.decode(z, pos_edge_index,
14             neg_edge_index)
15             link_probs = link_logits.sigmoid()
16             link_labels = get_link_labels(pos_edge_index,
17             neg_edge_index)
18             results.append(roc_auc_score(link_labels.cpu(),
19             link_probs.cpu()))
20     return results
```

注：在验证与测试阶段，我们也应该只见训练集，对验证集与测试集都是不可见的。所以在验证与测试阶段，我们依然只用训练集正样本边。

运行完整的训练、验证与测试

```
1  def main():
```

```

2     device = torch.device('cuda' if
torch.cuda.is_available() else 'cpu')
3
4     dataset = 'Cora'
5     path =
osp.join(osp.dirname(osp.realpath(__file__)), '..',
'data', dataset)
6     dataset = Planetoid(path, dataset,
transform=T.NormalizeFeatures())
7     data = dataset[0]
8     ground_truth_edge_index =
data.edge_index.to(device)
9     data.train_mask = data.val_mask = data.test_mask =
data.y = None
10    data = train_test_split_edges(data)
11    data = data.to(device)
12
13    model = Net(dataset.num_features, 64).to(device)
14    optimizer =
torch.optim.Adam(params=model.parameters(), lr=0.01)
15
16    best_val_auc = test_auc = 0
17    for epoch in range(1, 101):
18        loss = train(data, model, optimizer)
19        val_auc, tmp_test_auc = test(data, model)
20        if val_auc > best_val_auc:
21            best_val_auc = val_auc
22            test_auc = tmp_test_auc
23        print(f'Epoch: {epoch:03d}, Loss: {loss:.4f},
val: {val_auc:.4f}, '
24              f'Test: {test_auc:.4f}')
25
26    z = model.encode(data.x, data.train_pos_edge_index)
27    final_edge_index = model.decode_all(z)
28
29
30 if __name__ == "__main__":
31     main()
32

```

完整的代码可见于 `codes/edge_prediction.py`。

结语

在完整的第6节内容中，我们学习了

- PyG中规定的使用数据的一般过程；
- `InMemoryDataset` 基类；
- 一个简化的 `InMemory` 数据集类；
- 一个 `InMemory` 数据集类实例，以及使用该数据集类时会发生的一些过程；
- 节点预测任务实践；
- 边预测任务实践。

我们需要重点关注 `InMemory` 数据集类的运行流程与其四个方法的定义规范，同时我们还应该重点关注边预测任务中的数据集划分，训练集负样本采样，以及训练、验证与测试三个阶段使用的边。

作业

- 实践问题一：尝试使用PyG中的不同的网络层去代替 `GCNConv`，以及不同的层数和不同的 `out_channels`，来实现节点分类任务。
- 实践问题二：在边预测任务中，尝试用 `torch_geometric.nn.Sequential` 容器构造图神经网络。
- 思考问题三：如下方代码所示，我们以 `data.train_pos_edge_index` 为实际参数来进行训练集负样本采样，但这样采样得到的负样本可能包含一些验证集的正样本与测试集的正样本，即可能将真实的正样本标记为负样本，由此会产生冲突。但我们还是这么做，这是为什么？

```
1 neg_edge_index = negative_sampling(  
2     edge_index=data.train_pos_edge_index,  
3     num_nodes=data.num_nodes,  
4  
    num_neg_samples=data.train_pos_edge_index.size(1))
```

参考资料

- `Sequential` 官网文档：[torch_geometric.nn.Sequential](#)
- 边预测任务实践中的代码来源于[link_pred.py](#)

