

金融异常检测任务 - 程序报告

吴天宇 12334125

1 实验概要

1.1 实验内容

本实验旨在基于 DGraph-Fin 数据集（包含用户之间的社交网络关系和节点特征），利用图神经网络（Graph Neural Networks）和多层感知机（MLP）在金融领域进行异常检测，识别欺诈用户。

实验主要包括以下内容：

- 使用 PyTorch 和 PyTorch Geometric 进行图数据的加载和预处理。
- 定义并训练多层感知机（MLP）模型和 GraphSAGE 模型。
- 评估模型在节点分类任务中的性能，主要使用 AUC（Area Under the Curve）作为评估指标。
- 分析模型的训练过程和结果。

1.2 实验结果概要

在本实验中，我们分别训练了 MLP 模型和 GraphSAGE 模型。通过对比，我们发现：

- MLP 模型：只利用节点的特征信息，未考虑图结构，训练速度较快，但在验证集上的 AUC 表现有限。
- GraphSAGE 模型：结合了节点的特征和邻居信息，通过图卷积捕捉节点之间的关系，在验证集上取得了更高的 AUC。

最终，GraphSAGE 模型在验证集上取得了更优的性能，证明了利用图结构信息对于金融异常检测任务的重要性。

代码和报告开源于 https://github.com/Wuty-zju/zju_ai_sys

2 多层感知机 MLP 模型

本小节介绍了多层感知机（MLP）模型的构建和训练过程。导入了必要的库，包括 PyTorch、PyTorch Geometric 及自定义工具，并设置了计算设备（GPU 或 CPU）。加载并预处理数据集时，定义了数据路径，使用 DGraphFin 类加载 'DGraph' 数据集，对节点特征进行了标准化处理，并划分为训练、验证和测试集。定义了继承自 nn.Module 的 MLP 类，结构包含多个线性层、可选的批归一化层、ReLU 激活函数和 Dropout 层，最终通过 log_softmax 获取对数概率分布。在训练和评估函数部分，设置了训练超参数，采用 Adam 优化器和 AUC 评估器，明确了训练迭代和性能测试的方法。通过 train_model 函数，模型在指定轮数内进行训练，定期记录和输出损失值及 AUC 指标，并根据验证集的表现保存最佳模型。此外，提供了加载最佳模型权重和预测特定节点标签概率的功能，确保模型在实际应用中的可用性和准确性。

2.1 导入必要的库

```
In [ ]: import torch
import torch.nn.functional as F
import torch.nn as nn
import torch_geometric.transforms as T
from utils import DGraphFin
from utils.evaluator import Evaluator
import os

device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
```

2.2 数据加载和预处理

- 设置路径和检查目录：** 定义数据集保存路径 path 和模型保存路径 save_dir，如果结果目录不存在则创建。
- 加载数据集：** 指定数据集名称 'DGraph'，使用 DGraphFin 类加载数据集，应用稀疏张量转换 T.ToSparseTensor()，获取数据对象 data。
- 数据预处理：** 提取节点特征 x，对其进行标准化处理（减去均值并除以标准差），使特征具有零均值和单位方差，然后将处理后的特征赋值回数据对象。
- 划分数据集：** 根据数据的掩码属性，创建字典 split_idx，包含训练集、验证集和测试集的索引，获取训练集索引

```
train_idx。
```

- **数据迁移到计算设备：** 使用 `data = data.to(device)` 将数据对象移动到指定的计算设备（GPU 或 CPU）上，为后续的训练做准备。

```
In [ ]: # ===== 数据加载和预处理 =====
# 路径和参数设置
path = './datasets/632d74d4e2843a53167ee9a1-momodel/' # 数据保存路径
save_dir = './results/' # 模型保存路径
if not os.path.exists(save_dir):
    os.makedirs(save_dir)
dataset_name = 'DGraph' # 数据集名称

dataset = DGraphFin(root=path, name=dataset_name, transform=T.ToSparseTensor())
data = dataset[0]

# 数据预处理
x = data.x
x = (x - x.mean(0)) / x.std(0) # 标准化节点特征
data.x = x

# 划分训练集、验证集和测试集
split_idx = {
    'train': data.train_mask,
    'valid': data.valid_mask,
    'test': data.test_mask
}
train_idx = split_idx['train']

# 将数据移动到设备上 (GPU 或 CPU)
data = data.to(device)
```

2.3 定义模型

- **定义 MLP 模型类：** 创建名为 `MLP` 的类，继承自 `nn.Module`，用于构建多层感知机模型。
- **初始化模型参数：** 在 `__init__` 方法中，接收输入维度 `in_channels`、隐藏层维度 `hidden_channels`、输出维度 `out_channels`、层数 `num_layers`、Dropout 概率 `dropout` 和是否使用批归一化 `batchnorm` 等参数。初始化线性层列表 `self.lins`，添加从输入层到隐藏层的线性映射。根据指定的层数，循环添加隐藏层的线性映射和批归一化层（如果使用批归一化）。最后添加输出层的线性映射。设置 Dropout 概率。
- **重置模型参数：** 定义 `reset_parameters` 方法，重置所有线性层和批归一化层的参数，以初始化模型。
- **定义前向传播过程：** 在 `forward` 方法中，实现模型的前向传播过程。依次对输入数据应用线性变换、批归一化（如果使用）、ReLU 激活函数和 Dropout 操作。最后一层不经过激活和 Dropout，直接应用线性变换。使用 `F.log_softmax` 对模型输出进行处理，得到对数概率分布。

```
In [ ]: # ===== 定义模型 =====
class MLP(nn.Module):
    def __init__(self, in_channels, hidden_channels, out_channels, num_layers, dropout, batchnorm=True):
        super(MLP, self).__init__()
        # 定义多层感知机结构
        self.lins = nn.ModuleList([nn.Linear(in_channels, hidden_channels)])
        self.bns = nn.ModuleList([nn.BatchNorm1d(hidden_channels)] if batchnorm else None)
        for _ in range(num_layers - 2):
            self.lins.append(nn.Linear(hidden_channels, hidden_channels))
            if batchnorm:
                self.bns.append(nn.BatchNorm1d(hidden_channels))
        self.lins.append(nn.Linear(hidden_channels, out_channels))
        self.dropout = dropout

    def reset_parameters(self):
        # 重置模型参数
        for lin in self.lins:
            lin.reset_parameters()
        if self.bns:
            for bn in self.bns:
                bn.reset_parameters()

    def forward(self, x):
        # 模型前向传播
        for i, lin in enumerate(self.lins[:-1]):
            x = lin(x)
            if self.bns:
                x = self.bns[i](x)
            x = F.relu(x)
```

```
x = F.dropout(x, p=self.dropout, training=self.training)
x = self.lins[-1](x)
return F.log_softmax(x, dim=-1)
```

2.4 训练和评估函数

- **设置训练超参数和模型实例化**：定义模型的层数 `num_layers`、隐藏层维度 `hidden_channels` 等参数，创建包含这些参数的字典 `mlp_parameters`。然后，获取输入和输出维度，实例化 MLP 模型并将其移动到计算设备上。
- **定义优化器和损失函数**：使用 `torch.optim.Adam` 优化器，设置学习率 `lr` 和权重衰减 `weight_decay`。
- **定义评估器**：使用 AUC（曲线下面积）作为评估指标，创建评估器实例 `Evaluator('auc')`。
- **定义训练函数 `train`**：该函数对模型进行一次训练迭代，包括模型切换到训练模式、梯度清零、前向传播、计算损失（使用负对数似然损失 `F.nll_loss`）、反向传播和参数更新。返回当前的损失值。
- **定义测试函数 `test`**：该函数在训练集和验证集上评估模型性能。模型切换到评估模式，禁止梯度计算。对指定的数据集进行前向传播，计算预测概率和损失值。使用评估器计算 AUC 分数，并返回评估结果、损失值和预测结果。

```
In [ ]: # ===== 训练和评估函数 =====
# 训练超参数设置
num_layers = 5
hidden_channels = 128

mlp_parameters = {'num_layers': num_layers, 'hidden_channels': hidden_channels, 'dropout': 0.5, 'batchnorm':
in_channels, out_channels = data.x.size(-1), 2
model = MLP(in_channels, **mlp_parameters, out_channels=out_channels).to(device)

# 优化器和损失函数
optimizer = torch.optim.Adam(model.parameters(), lr=0.001, weight_decay=5e-4)

# 评估器
evaluator = Evaluator('auc')

def train(model, data, train_idx, optimizer):
    """
    训练模型
    :param model: 模型对象
    :param data: 数据对象
    :param train_idx: 训练集索引
    :param optimizer: 优化器
    :return: 损失值
    """
    model.train()
    optimizer.zero_grad()
    out = model(data.x[train_idx])
    loss = F.nll_loss(out, data.y[train_idx].squeeze().long())
    loss.backward()
    optimizer.step()
    return loss.item()

def test(model, data, split_idx, evaluator):
    """
    测试模型性能
    :param model: 模型对象
    :param data: 数据对象
    :param split_idx: 数据集划分字典
    :param evaluator: 评估器
    :return: 评估结果、损失和预测值
    """
    model.eval()
    with torch.no_grad():
        losses, eval_results = {}, {}
        for key in ['train', 'valid']:
            node_id = split_idx[key]
            out = model(data.x[node_id])
            y_pred = out.exp()
            losses[key] = F.nll_loss(out, data.y[node_id].squeeze().long()).item()
            eval_results[key] = evaluator.eval(data.y[node_id].squeeze().long(), y_pred)['auc']
    return eval_results, losses, y_pred
```

2.5 训练模型

- **定义训练模型函数 `train_model`**：创建一个函数，负责执行模型的训练过程并保存最佳模型。
- **设置最佳指标和训练集索引**：初始化 `best_valid_auc` 为0和 `min_valid_loss` 为无限大，获取训练集索引 `train_idx`。

- **训练循环**：在指定的训练轮数 `epochs` 内进行迭代，每轮执行训练和测试，更新最佳模型并根据验证损失保存模型。
- **打印训练状态**：每隔 `log_steps` 轮输出当前轮次的损失值、训练集AUC和验证集AUC。
- **启动训练**：调用 `train_model` 函数，传入模型、数据、划分索引、优化器、评估器和保存目录，开始模型训练。

```
In [ ]: # ===== 训练模型 =====
def train_model(model, data, split_idx, optimizer, evaluator, save_dir, epochs=1000, log_steps=10):
    """
    执行模型训练并保存最佳模型。
    :param model: 模型对象
    :param data: 数据对象
    :param split_idx: 数据集划分字典
    :param optimizer: 优化器
    :param evaluator: 评估器
    :param save_dir: 模型保存路径
    :param epochs: 训练轮数
    :param log_steps: 日志记录频率
    """

    best_valid_auc, min_valid_loss = 0, float('inf')
    train_idx = split_idx['train']

    for epoch in range(1, epochs + 1):
        loss = train(model, data, train_idx, optimizer)
        eval_results, losses, _ = test(model, data, split_idx, evaluator)
        train_auc, valid_auc = eval_results['train'], eval_results['valid']
        train_loss, valid_loss = losses['train'], losses['valid']

        # 保存最优模型
        if valid_loss < min_valid_loss:
            min_valid_loss = valid_loss
            torch.save(model.state_dict(), os.path.join(save_dir, f'best_mlp_model_layers{num_layers}_hidden{num_hidden}.pt'))

        if epoch % 10 == 0:
            print(f'第 {epoch:04d} 轮, 损失值: {loss:.4f}, 训练集 AUC: {train_auc * 100:.2f}%, 验证集 AUC: {valid_auc * 100:.2f}%')

    train_model(model, data, split_idx, optimizer, evaluator, save_dir)
```

2.6 保存并加载最佳模型

- **定义加载最佳模型函数**：创建 `load_best_model` 函数，用于从指定的保存目录加载最佳模型权重。
- **加载模型权重**：使用 `model.load_state_dict` 从文件 `'best_mlp_model.pt'` 中加载权重。
- **返回加载后的模型**：函数返回加载了权重的模型对象。

```
In [ ]: # ===== 保存并加载最佳模型 =====
def load_best_model(model, save_dir):
    """
    加载最佳模型权重。
    :param model: 模型对象
    :param save_dir: 模型保存路径
    :return: 加载权重后的模型
    """

    model.load_state_dict(torch.load(os.path.join(save_dir, 'best_mlp_model.pt')))
    return model
```

2.7 测试函数

- **定义预测函数 `predict`**：创建 `predict` 函数，用于预测指定节点的标签概率。
- **设置模型为评估模式**：调用 `model.eval()`，将模型切换到评估状态，关闭诸如 dropout 的训练特性。
- **禁用梯度计算**：使用 `with torch.no_grad()`，在预测过程中不计算梯度，以节省内存和计算资源。
- **执行前向传播**：对指定节点的数据进行前向传播，获取模型输出。
- **转换输出为概率**：使用 `out.exp()` 将对数概率转换为实际概率分布。
- **返回预测概率**：返回节点的预测概率 `y_pred`。

```
In [ ]: # ===== 测试函数 =====
def predict(model, data, node_id):
    """
    预测指定节点的标签概率。
    :param model: 训练好的模型
    """
```



```

:param data: 数据对象
:param node_id: 节点索引
:return: 节点的预测概率
"""
model.eval()
with torch.no_grad():
    out = model(data.x[node_id].unsqueeze(0))
    y_pred = out.exp()
return y_pred

```

3 图神经网络 GraphSAGE 模型

本小节实现图神经网络 GraphSAGE 模型。首先加载并预处理图数据，包括标准化节点特征和将有向图转换为无向图。接着定义了一个包含三个 SAGEConv 层和残差连接的 GraphSAGE 模型，并配置了训练的超参数和优化器。通过定义训练和测试函数，模型在训练集和验证集上进行迭代训练与评估，保存验证集 AUC 最佳的模型状态，并在前向传播后保存预测结果。最后，提供了测试和预测函数，用于运行模型的前向传播、保存预测结果以及返回指定节点的类别概率。

3.1 导入必要的库

```
In [ ]: import torch
import torch.nn as nn
import torch.nn.functional as F
from torch_geometric.nn import SAGEConv
from torch_geometric.data import Data
import torch_geometric.transforms as T
from utils import DGraphFin
from utils.evaluator import Evaluator
import numpy as np
import os

device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
```

3.2 数据加载和预处理

- **设置数据路径和模型保存路径：**定义数据集的存储路径 `path` 以及模型保存目录 `save_dir`，并检查保存目录是否存在，若不存在则创建。
- **加载数据集：**使用 `DGraphFin` 类加载名为 `DGraph` 的数据集，并将有向图转换为无向图。
- **数据预处理：**
 - **标准化节点特征：**对节点特征 `x` 进行均值为 0、标准差为 1 的标准化处理。
 - **调整标签维度：**如果标签 `y` 的维度为 2，则压缩为一维。
- **划分数据集：**根据 `train_mask`、`valid_mask` 和 `test_mask` 划分训练集、验证集和测试集的索引。
- **移动数据到设备：**将数据移动到可用的计算设备（GPU 或 CPU）上，以加速计算。
- **转换邻接矩阵格式：**将稀疏的邻接矩阵 `adj_t` 转换为适用于 `SAGEConv` 的 `edge_index` 格式，便于后续的图卷积操作。

```
In [ ]: # ===== 数据加载和预处理 =====
# 数据路径设置
path = './datasets/632d74d4e2843a53167ee9a1-momodel/' # 数据保存路径
save_dir = './results/' # 模型保存路径
if not os.path.exists(save_dir):
    os.makedirs(save_dir)
dataset_name = 'DGraph' # 数据集名称

# 加载数据集
dataset = DGraphFin(root=path, name=dataset_name, transform=T.ToSparseTensor())
nlabels = 2 # 仅需预测类别 0 和类别 1
data = dataset[0]
data.adj_t = data.adj_t.to_symmetric() # 将有向图转换为无向图

# 数据预处理
x = data.x
x = (x - x.mean(0)) / x.std(0) # 标准化节点特征
data.x = x
if data.y.dim() == 2:
    data.y = data.y.squeeze(1) # 如果标签维度为 2，则压缩为 1 维

# 划分训练集、验证集和测试集
split_idx = {
    'train': data.train_mask,
```

```

        'valid': data.valid_mask,
        'test': data.test_mask
    }
    train_idx = split_idx['train']

    # 将数据移动到设备上 (GPU 或 CPU)
    data = data.to(device)

    # 将稀疏邻接矩阵 adj_t 转换为 edge_index (适用于 SAGEConv)
    row, col, _ = data.adj_t.coo() # 获取 COO 格式的行、列索引
    data.edge_index = torch.stack([row, col], dim=0) # 构建 edge_index 矩阵, 形状为 [2, num_edges]

```

3.3 定义模型

- 模型结构：
 - 创建 `GraphSAGE` 类，继承自 `nn.Module`，用于构建图神经网络模型。
 - 定义三个 `SAGEConv` 层：
 - 第一层将输入特征维度转换为隐藏层维度。
 - 第二层保持隐藏层维度不变。
 - 第三层将隐藏层维度转换为输出类别数。
 - 定义用于残差连接的线性层：
 - `res1`：用于在输入特征维度与隐藏层维度不同时进行调整。
 - `res2`：用于在两个隐藏层之间进行连接。
- 前向传播与参数重置：
 - 前向传播方法 `forward`：
 - 第一层卷积与残差连接：应用 `conv1` 卷积层和 ReLU 激活函数，调整残差连接的维度并与卷积输出相加。
 - 第二层卷积与残差连接：应用 `conv2` 卷积层和 ReLU 激活函数，使用 `res2` 进行残差连接并与卷积输出相加。
 - 第三层卷积与输出：应用 `conv3` 卷积层，使用 Log Softmax 函数获取类别概率。
 - 重置参数方法 `reset_parameters`：重置所有 `SAGEConv` 层和残差连接的线性层的参数，确保模型初始化状态。
- 模型实例化：
 - 根据输入特征维度、隐藏层维度和输出类别数初始化 `GraphSAGE` 模型。
 - 将模型移动到指定的计算设备（GPU 或 CPU）上，以加速计算。

```

In [ ]: # ===== 定义模型 =====
class GraphSAGE(nn.Module):
    def __init__(self, in_channels, hidden_channels, out_channels):
        super(GraphSAGE, self).__init__()
        # 定义三个 SAGEConv 层
        self.conv1 = SAGEConv(in_channels, hidden_channels)
        self.conv2 = SAGEConv(hidden_channels, hidden_channels)
        self.conv3 = SAGEConv(hidden_channels, out_channels)

        # 定义用于残差连接的线性层
        self.res1 = nn.Linear(in_channels, hidden_channels) if in_channels != hidden_channels else None
        self.res2 = nn.Linear(hidden_channels, hidden_channels)

    def reset_parameters(self):
        # 重置模型参数
        self.conv1.reset_parameters()
        self.conv2.reset_parameters()
        self.conv3.reset_parameters()
        if self.res1:
            self.res1.reset_parameters()
        self.res2.reset_parameters()

    def forward(self, x, edge_index):
        # 第一层卷积 + 残差连接
        identity = x # 保存输入以用于残差连接
        x = F.relu(self.conv1(x, edge_index)) # 图卷积和激活函数
        if self.res1:

```

```

        identity = self.res1(identity) # 如果维度不同, 调整维度
        x1 = x + identity # 残差连接

        # 第二层卷积 + 残差连接
        identity = x1
        x = F.relu(self.conv2(x1, edge_index))
        x2 = x + self.res2(identity) # 残差连接

        # 第三层卷积 (输出层)
        x3 = self.conv3(x2, edge_index)

        # 使用 Log Softmax 获取类别概率
        return F.log_softmax(x3, dim=-1)

# 实例化模型并移动到设备上
in_channels = data.x.size(-1) # 输入特征维度
hidden_channels = 2 # 隐藏层维度
out_channels = nlabels # 输出类别数
model = GraphSAGE(
    in_channels=in_channels,
    hidden_channels=hidden_channels,
    out_channels=out_channels
).to(device)

```

3.4 训练和评估函数

- 训练超参数设置：
 - 设置训练轮数（epochs）、学习率（lr）和权重衰减（weight_decay）。
- 优化器和损失函数：
 - 使用 Adam 优化器，并设置学习率和权重衰减。
 - 选择负对数似然损失函数（NLLLoss）作为损失函数。
- 评估器：
 - 使用 AUC 作为评估指标，初始化评估器。
- 训练函数 **train**：
 - 将模型设置为训练模式。
 - 清空梯度，进行前向传播计算输出。
 - 计算训练集上的损失，并进行反向传播更新参数。
 - 返回当前训练轮次的损失值。
- 测试函数 **test**：
 - 将模型设置为评估模式。
 - 禁用梯度计算，进行前向传播计算输出。
 - 将 Log Softmax 输出转换为概率。
 - 计算训练集和验证集上的损失，并评估 AUC 指标。
 - 返回评估结果和损失值。

```

In [ ]: # ===== 训练和评估函数 =====
# 训练超参数设置
epochs = 20 # 训练轮数
lr = 0.005 # 学习率
weight_decay = 2e-4 # 权重衰减 (L2 正则化系数)

# 优化器和损失函数
optimizer = torch.optim.Adam(model.parameters(), lr=lr, weight_decay=weight_decay)

# 评估器
eval_metric = 'auc' # 使用 AUC 作为评估指标
evaluator = Evaluator(eval_metric)

# 定义训练函数
def train(model, data, train_idx, optimizer):
    model.train() # 设置模型为训练模式
    optimizer.zero_grad() # 清空梯度

```

```

out = model(data.x, data.edge_index) # 前向传播
loss = F.nll_loss(out[train_idx], data.y[train_idx]) # 计算损失 (负对数似然损失)
loss.backward() # 反向传播
optimizer.step() # 更新参数
return loss.item() # 返回损失值

# 定义测试函数
def test(model, data, split_idx, evaluator):
    model.eval() # 设置模型为评估模式
    with torch.no_grad():
        out = model(data.x, data.edge_index) # 前向传播
        y_pred = out.exp() # 将 Log Softmax 输出转换为概率
        eval_results = {}
        losses = {}
        for key in ['train', 'valid']:
            node_id = split_idx[key]
            losses[key] = F.nll_loss(out[node_id], data.y[node_id]).item() # 计算损失
            # 计算评估指标 (AUC)
            eval_results[key] = evaluator.eval(data.y[node_id], y_pred[node_id])[eval_metric]
    return eval_results, losses # 返回评估结果和损失

```

3.5 训练、保存并加载最佳模型

- 初始化：
 - 初始化最佳验证集 AUC (`best_valid_auc`) 和最佳模型状态 (`best_model_state`)。
- 训练与评估循环：
 - 在每个训练轮次中，调用 `train` 函数进行模型训练，并调用 `test` 函数在训练集和验证集上进行评估。
 - 记录训练集和验证集的 AUC 指标。
- 保存最佳模型：
 - 如果当前验证集 AUC 优于之前的最佳 AUC，则更新最佳 AUC 并保存当前模型状态。
 - 将最佳模型状态保存到指定文件中。
- 打印训练进度：
 - 每10个轮次打印一次当前轮次的损失值、训练集 AUC 和验证集 AUC。
- 加载最佳模型：
 - 训练完成后，打印最佳验证集 AUC，并从保存的文件中加载最佳模型状态。

```

In [ ]: # ===== 训练、保存并加载最佳模型 =====
best_valid_auc = 0 # 初始化最佳验证集 AUC
best_model_state = None # 用于保存最佳模型状态

for epoch in range(1, epochs + 1):
    loss = train(model, data, train_idx, optimizer) # 训练一步
    eval_results, losses = test(model, data, split_idx, evaluator) # 在训练集和验证集上测试
    train_auc = eval_results['train']
    valid_auc = eval_results['valid']

    if valid_auc > best_valid_auc:
        best_valid_auc = valid_auc
        best_model_state = model.state_dict() # 保存当前最佳模型状态
        # 保存最佳模型
        model_filename = f'best_sage_model_conv3_hidden{hidden_channels}_lr{lr}_wd{weight_decay}.pt'
        torch.save(best_model_state, os.path.join(save_dir, model_filename))

    if epoch % 10 == 0:
        print(f'第 {epoch:04d} 轮, 损失值: {loss:.4f}, 训练集 AUC: {train_auc * 100:.2f}%, 验证集 AUC: {valid_auc * 100:.2f}%')

print("训练完成。")
print(f"最佳验证集 AUC: {best_valid_auc * 100:.2f}%")
print(f"最佳模型已保存至 {os.path.join(save_dir, model_filename)}")

model.load_state_dict(torch.load(os.path.join(save_dir, model_filename), map_location=device))

```

3.6 测试并保存预测结果函数

- 定义函数 `test_and_save_predictions`：该函数用于运行模型的前向传播，并保存所有节点的预测结果。
- 设置模型为评估模式：调用 `model.eval()` 将模型切换到评估模式，关闭诸如 dropout 的训练特性。

- **禁用梯度计算**: 使用 `with torch.no_grad()` 在预测过程中不计算梯度, 以节省内存和计算资源。
- **执行前向传播**: 对所有节点的数据进行前向传播, 获取模型输出 `out` 。
- **转换输出为概率**: 使用 `out.exp()` 将对数概率转换为实际概率分布 `y_pred` 。
- **保存预测结果**: 使用 `torch.save` 将预测结果保存到指定路径 `save_path` , 并打印保存成功的消息。
- **运行并保存预测结果**: 调用 `test_and_save_predictions` 函数, 运行模型并保存预测结果到指定文件路径。

```
In [ ]: # ===== 测试并保存预测结果函数 =====
def test_and_save_predictions(model, data, save_path):
    """
    运行模型的前向传播, 并保存所有节点的预测结果
    :param model: 训练好的模型
    :param data: 包含节点特征和边的图数据
    :param save_path: 保存预测结果的文件路径
    """
    model.eval() # 设置模型为评估模式
    with torch.no_grad():
        # 对所有节点进行前向传播
        out = model(data.x, data.edge_index)
        y_pred = out.exp() # 将 Log Softmax 输出转换为概率

    # 保存预测结果
    torch.save(y_pred.cpu(), save_path)
    print(f"预测结果已保存至 {save_path}")

# 运行模型并保存预测结果
predictions_save_path = os.path.join(
    save_dir,
    f'best_sage_model_conv3_hidden{hidden_channels}_lr{lr}_wd{weight_decay}_predictions.pt'
)
test_and_save_predictions(model, data, predictions_save_path)
```

3.7 测试-预测函数

- 定义函数 `predict` , 用于加载模型并在指定节点上进行预测。
- 使用 `torch.load` 加载保存的最佳模型权重。
- 对指定节点的数据进行前向传播, 获取模型输出 `out` 。
- 使用 `out[node_id].exp()` 将对数概率转换为实际概率分布 `y_pred` 。
- 返回指定节点的类别 0 和类别 1 的预测概率 `y_pred` 。

```
In [ ]: # ===== 测试-预测函数 =====
def predict(data, node_id):
    """
    加载模型并在 MoAI 平台进行预测
    :param data: 数据对象, 包含 x 和 edge_index 等属性
    :param node_id: int, 需要进行预测的节点索引
    :return: tensor, 类别 0 和类别 1 的概率
    """
    out = model
    y_pred = out[node_id].exp() # 获取指定节点的预测概率, 并增加一个维度

    return y_pred # 返回预测概率

model = torch.load('./results/best_sage_model_conv3_hidden128_lr0.002_wd0.0002_predictions.pt')
```