Московский государственный технический университет им. Н.Э. Баумана
Факультет «Информатика и системы управления»
Кафедра «Системы обработки информации и управления»

# Лабораторная работа №6

## «Разработка системы предсказаний поведения на основании графовых моделей»

по дисциплине

«Методы машинного обучения»

**ИСПОЛНИТЕЛЬ:**

Крюков Г.М.
Группа ИУ5-21М

_____

"__"_____2022 г.

Москва, 2022

**Цель:** Обучение работе с предварительной обработкой графовых типов данных и обучением нейронных сетей на графовых данных.

## Задание:

Подготовить датасет графовых данных

Подобрать модель и гиперпараметры обучения для получения качества AUC > 0.65

## Набор данных:

В качестве базы данных предлагаем использовать датасет о покупках пользователей в одном магазине товаров RecSys Challenge 2015 (https://www.kaggle.com/datasets/chadgostopp/recsys-challenge-2015).

# Текст программы

## ▾ Лабораторная работа №6:

### "Разработка системы предсказания поведения на основании графовых моделей"

*Цель*: обучение работе с графовым типом данных и графовыми нейронными сетями.

*Задача*: подготовить графовый датасет из базы данных о покупках и построить модель предсказания совершения покупки.

## Графовые нейронные сети

**Графовые нейронные сети** - тип нейронной сети, которая напрямую работает со структурой графа. Типичным применениями GNN являются:

- Классификация узлов;
- Предсказание связей;
- Графовая классификация;
- Распознавание движений;
- Рекомендательные системы.

В данной лабораторной работе будет происходить работа над **графовыми сверточными сетями**. Отличаются они от сверточных нейронных сетей нефиксированной структурой, функция свертки не является .

Подробнее можно прочитать тут: https://towardsdatascience.com/understanding-graph-convolutional-networks-for-node-classification-a2bfdb7aba7b

Тут можно почитать современные подходы к использованию графовых сверточных сетей
https://paperswithcode.com/method/gcn

## Датасет

В качестве базы данных предлагаем использовать датасет о покупках пользователей в одном магазине товаров RecSys Challenge 2015 (https://www.kaggle.com/datasets/chadgostopp/recsys-challenge-2015).

Скачать датасет можно отсюда: https://drive.google.com/drive/folders/1gtAeXPTj-c0RwVOKreMrZ3bfSmCwl2y-?usp=sharing (lite-версия является облегченной версией исходного датасета, рекомендуем использовать её)

Также рекомендуем загружать данные в виде архива и распаковывать через пакет zipfile или/и скачивать датасет в собственный Google Drive и примонтировать его в колаб.

## ▾ Установка библиотек, выгрузка исходных датасетов

```
# Slow method of installing pytorch geometric
# !pip install torch_geometric
# !pip install torch_sparse
# !pip install torch_scatter

# Install pytorch geometric
!pip install torch-sparse -f https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
!pip install torch-cluster -f https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
!pip install torch-spline-conv -f https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
!pip install torch-geometric -f https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
!pip install torch-scatter==2.0.9 -f https://data.pyg.org/whl/torch-1.11.0%2Bcu113.html
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Looking in links: https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
Collecting torch-sparse
  Downloading https://data.pyg.org/whl/torch-1.11.0%2Bcu113/torch_sparse-0.6.13-cp37m-linux_x86_64.whl (3.5 MB)
     |████████████████████████████████| 3.5 MB 1.3 MB/s
Requirement already satisfied: scipy in /usr/local/lib/python3.7/dist-packages (from torch-sparse) (1.4.1)
Requirement already satisfied: numpy>=1.13.3 in /usr/local/lib/python3.7/dist-packages (from scipy->torch-sparse) (1
Installing collected packages: torch-sparse
Successfully installed torch-sparse-0.6.13
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Looking in links: https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
Collecting torch-cluster
  Downloading https://data.pyg.org/whl/torch-1.11.0%2Bcu113/torch_cluster-1.6.0-cp37m-linux_x86_64.whl (2.5 MB)
     |████████████████████████████████| 2.5 MB 26.3 MB/s
Installing collected packages: torch-cluster
```

```
Successfully installed torch-cluster-1.6.0
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Looking in links: https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
Collecting torch-spline-conv
  Downloading https://data.pyg.org/whl/torch-1.11.0%2Bcu113/torch_spline_conv-1.2.1-cp37-cp37m-linux_x86_64.whl (750
|███████████████████████████████| 750 kB 23.7 MB/s
Installing collected packages: torch-spline-conv
Successfully installed torch-spline-conv-1.2.1
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Looking in links: https://pytorch-geometric.com/whl/torch-1.11.0%2Bcu113.html
Collecting torch-geometric
  Downloading torch_geometric-2.0.4.tar.gz (407 kB)
|███████████████████████████████| 407 kB 25.7 MB/s
Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (4.64.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (1.21.6)
Requirement already satisfied: scipy in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (1.4.1)
Requirement already satisfied: pandas in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (1.3.5)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (2.11.3)
Requirement already satisfied: requests in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (2.23.0)
Requirement already satisfied: pyparsing in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (3.0.9)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.7/dist-packages (from torch-geometric) (1.0.2)
Requirement already satisfied: MarkupSafe>=0.23 in /usr/local/lib/python3.7/dist-packages (from jinja2->torch-geomet
Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/dist-packages (from pandas->torch-geometric)
Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-packages (from pandas->torch-
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-packages (from python-dateutil>=2.7.3->pand
Requirement already satisfied: urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1 in /usr/local/lib/python3.7/dist-packages (fr
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from requests->torch-ge
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests->torch-geo
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from requests->torch-geometri
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (from scikit-learn->torch-geom
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from scikit-learn->to
Building wheels for collected packages: torch-geometric
  Building wheel for torch-geometric (setup.py) ... done
  Created wheel for torch-geometric: filename=torch_geometric-2.0.4-py3-none-any.whl size=616603 sha256=0eaf2157f99f
  Stored in directory: /root/.cache/pip/wheels/18/a6/a4/ca18c3051fcead866fe7b85700ee2240d883562a1bc70ce421
Successfully built torch-geometric
Installing collected packages: torch-geometric
Successfully installed torch-geometric-2.0.4
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Looking in links: https://data.pyg.org/whl/torch-1.11.0%2Bcu113.html
Collecting torch-scatter==2.0.9
```

```python
import numpy as np
import pandas as pd
import pickle
import csv
import os

from sklearn.preprocessing import LabelEncoder

import torch

# PyG - PyTorch Geometric
from torch_geometric.data import Data, DataLoader, InMemoryDataset

from tqdm import tqdm


RANDOM_SEED = 17 #@param { type: "integer" }
BASE_DIR = '/content/' #@param { type: "string" }
np.random.seed(RANDOM_SEED)


# Check if CUDA is available for colab
torch.cuda.is_available

    <function torch.cuda.is_available>


# Unpack files from zip-file
import zipfile
with zipfile.ZipFile(BASE_DIR + 'yoochoose-data-lite.zip', 'r') as zip_ref:
    zip_ref.extractall(BASE_DIR)
```

RANDOM_SEED: 17

BASE_DIR: " /content/ "

▾ Анализ исходных данных

```python
# Read dataset of items in store
df = pd.read_csv(BASE_DIR + 'yoochoose-clicks-lite.dat')
```

```python
# df.columns = ['session_id', 'timestamp', 'item_id', 'category']
df.head()
```

```
/usr/local/lib/python3.7/dist-packages/IPython/core/interactiveshell.py:2882: DtypeWarning: Columns (3) have mixed type
  exec(code_obj, self.user_global_ns, self.user_ns)
```

| | session_id | timestamp | item_id | category |
|---|---|---|---|---|
| 0 | 9 | 2014-04-06T11:26:24.127Z | 214576500 | 0 |
| 1 | 9 | 2014-04-06T11:28:54.654Z | 214576500 | 0 |
| 2 | 9 | 2014-04-06T11:29:13.479Z | 214576500 | 0 |
| 3 | 19 | 2014-04-01T20:52:12.357Z | 214561790 | 0 |
| 4 | 19 | 2014-04-01T20:52:13.758Z | 214561790 | 0 |

```python
# Read dataset of purchases
buy_df = pd.read_csv(BASE_DIR + 'yoochoose-buys-lite.dat')
# buy_df.columns = ['session_id', 'timestamp', 'item_id', 'price', 'quantity']
buy_df.head()
```

| | session_id | timestamp | item_id | price | quantity |
|---|---|---|---|---|---|
| 0 | 420374 | 2014-04-06T18:44:58.314Z | 214537888 | 12462 | 1 |
| 1 | 420374 | 2014-04-06T18:44:58.325Z | 214537850 | 10471 | 1 |
| 2 | 489758 | 2014-04-06T09:59:52.422Z | 214826955 | 1360 | 2 |
| 3 | 489758 | 2014-04-06T09:59:52.476Z | 214826715 | 732 | 2 |
| 4 | 489758 | 2014-04-06T09:59:52.578Z | 214827026 | 1046 | 1 |

```python
# Filter out item session with length < 2
df['valid_session'] = df.session_id.map(df.groupby('session_id')['item_id'].size() > 2)
df = df.loc[df.valid_session].drop('valid_session',axis=1)
df.nunique()
```

```
session_id    1000000
timestamp     5557758
item_id         37644
category          275
dtype: int64
```

```python
# Randomly sample a couple of them
NUM_SESSIONS = 60000 #@param { type: "integer" }
sampled_session_id = np.random.choice(df.session_id.unique(), NUM_SESSIONS, replace=False)
df = df.loc[df.session_id.isin(sampled_session_id)]
df.nunique()
```

NUM_SESSIONS: 60000

```
session_id     60000
timestamp     334990
item_id        20043
category         103
dtype: int64
```

```python
# Average length of session
df.groupby('session_id')['item_id'].size().mean()
```

```
5.5834166666666665
```

```python
# Encode item and category id in item dataset so that ids will be in range (0,len(df.item.unique()))
item_encoder = LabelEncoder()
category_encoder = LabelEncoder()
df['item_id'] = item_encoder.fit_transform(df.item_id)
df['category']= category_encoder.fit_transform(df.category.apply(str))
df.head()
```

| | session_id | timestamp | item_id | category |
|---|---|---|---|---|
| **91** | 131 | 2014-04-03T04:46:08.891Z | 13649 | 0 |

```
# Encode item and category id in purchase dataset
buy_df = buy_df.loc[buy_df.session_id.isin(df.session_id)]
buy_df['item_id'] = item_encoder.transform(buy_df.item_id)
buy_df.head()
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:3: SettingWithCopyWar
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable
 This is separate from the ipykernel package so we can avoid doing imports until
```

| | session_id | timestamp | item_id | price | quantity |
|---|---|---|---|---|---|
| **5** | 70427 | 2014-04-02T15:54:07.144Z | 13729 | 3769 | 1 |
| **25** | 140964 | 2014-04-04T07:02:02.655Z | 10268 | 2408 | 1 |
| **62** | 489671 | 2014-04-03T15:48:37.392Z | 13710 | 4188 | 1 |
| **63** | 489671 | 2014-04-03T15:59:35.495Z | 13710 | 4188 | 1 |
| **64** | 489671 | 2014-04-03T16:00:06.917Z | 13710 | 4188 | 1 |

```
# Get item dictionary with grouping by session
buy_item_dict = dict(buy_df.groupby('session_id')['item_id'].apply(list))
buy_item_dict
```

```
{714: [16129, 16324, 16326, 3323],
 3517: [11939, 13381],
 4832: [12191, 12191, 12191],
 5002: [12217],
 5942: [16913, 14322, 14040, 14040, 14040, 14322, 16913],
 7054: [14021],
 7173: [13549],
 8628: [280, 281],
 9292: [13783, 4280],
```

```
 9702: [15763, 15756, 14250, 3096, 9004, 13708, 11207, 14092],
 10879: [2311, 2362],
 12017: [5219],
 12282: [13864,
  13866,
  13862,
  9463,
  13866,
  13864,
  13862,
  9463,
  13866,
  13864,
  13862,
  9463],
 13073: [12213, 12213],
 14227: [11645, 11646],
 14314: [12227, 13658, 13585, 13773],
 16053: [4365, 13874, 13872],
 16946: [8864],
 19576: [13862, 13862],
 21812: [1049],
 23934: [14142],
 24703: [14079, 14079],
 28668: [14065],
 35699: [1758],
 38117: [8306,
  12079,
  13907,
  13907,
  13618,
  13775,
  13775,
  13625,
  13625,
  12959,
  12960,
  14764],
 39247: [13862, 13866],
 40176: [4277],
 40209: [15763, 15762, 15756, 10293],
 40717: [15058],
```

```
        40827: [16129, 16129],
        42536: [14040, 14322],
        44097: [13520],
        44714: [1468, 13523, 14321, 13523],
        45836: [10705],
        46132: [14021, 2584],
```

▾ Сборка выборки для обучения

```python
# Transform df into tensor data
def transform_dataset(df, buy_item_dict):
    data_list = []

    # Group by session
    grouped = df.groupby('session_id')
    for session_id, group in tqdm(grouped):
        le = LabelEncoder()
        sess_item_id = le.fit_transform(group.item_id)
        group = group.reset_index(drop=True)
        group['sess_item_id'] = sess_item_id

        #get input features
        node_features = group.loc[group.session_id==session_id,
                            ['sess_item_id','item_id','category']].sort_values('sess_item_id')[['item_id','category'
        node_features = torch.LongTensor(node_features).unsqueeze(1)
        target_nodes = group.sess_item_id.values[1:]
        source_nodes = group.sess_item_id.values[:-1]

        edge_index = torch.tensor([source_nodes,
                            target_nodes], dtype=torch.long)
        x = node_features

        #get result
        if session_id in buy_item_dict:
            positive_indices = le.transform(buy_item_dict[session_id])
            label = np.zeros(len(node_features))
            label[positive_indices] = 1
```

```python
        else:
            label = [0] * len(node_features)

        y = torch.FloatTensor(label)

        data = Data(x=x, edge_index=edge_index, y=y)

        data_list.append(data)

    return data_list

# Pytorch class for creating datasets
class YooChooseDataset(InMemoryDataset):
    def __init__(self, root, transform=None, pre_transform=None):
        super(YooChooseDataset, self).__init__(root, transform, pre_transform)
        self.data, self.slices = torch.load(self.processed_paths[0])

    @property
    def raw_file_names(self):
        return []

    @property
    def processed_file_names(self):
        return [BASE_DIR+'yoochoose_click_binary_100000_sess.dataset']

    def download(self):
        pass

    def process(self):
        data_list = transform_dataset(df, buy_item_dict)

        data, slices = self.collate(data_list)
        torch.save((data, slices), self.processed_paths[0])


# Prepare dataset
dataset = YooChooseDataset('./')
```

```
Processing...
   0%|          | 0/60000 [00:00<?, ?it/s]/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:21: UserWarning:
 100%|██████████| 60000/60000 [03:23<00:00, 294.14it/s]
Done!
```

▾ Разделение выборки

```python
# train_test_split
dataset = dataset.shuffle()
one_tenth_length = int(len(dataset) * 0.1)
train_dataset = dataset[:one_tenth_length * 8]
val_dataset = dataset[one_tenth_length*8:one_tenth_length * 9]
test_dataset = dataset[one_tenth_length*9:]
len(train_dataset), len(val_dataset), len(test_dataset)
```

```
    (48000, 6000, 6000)
```

```python
# Load dataset into PyG loaders
batch_size= 512
train_loader = DataLoader(train_dataset, batch_size=batch_size)
val_loader = DataLoader(val_dataset, batch_size=batch_size)
test_loader = DataLoader(test_dataset, batch_size=batch_size)
```

```
    /usr/local/lib/python3.7/dist-packages/torch_geometric/deprecation.py:12: UserWarning: 'data.DataLoader' is deprecated,
      warnings.warn(out)
```

```python
# Load dataset into PyG loaders
num_items = df.item_id.max() +1
num_categories = df.category.max()+1
num_items , num_categories
```

```
    (20043, 102)
```

▾ Настройка модели для обучения

```python
embed_dim = 128
from torch_geometric.nn import GraphConv, TopKPooling, GatedGraphConv, SAGEConv, SGConv
from torch_geometric.nn import global_mean_pool as gap, global_max_pool as gmp
import torch.nn.functional as F

class Net(torch.nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        # Model Structure
        self.conv1 = GraphConv(embed_dim * 2, 128)
        self.pool1 = TopKPooling(128, ratio=0.9)
        self.conv2 = GraphConv(128, 128)
        self.pool2 = TopKPooling(128, ratio=0.9)
        self.conv3 = GraphConv(128, 128)
        self.pool3 = TopKPooling(128, ratio=0.9)
        self.item_embedding = torch.nn.Embedding(num_embeddings=num_items, embedding_dim=embed_dim)
        self.category_embedding = torch.nn.Embedding(num_embeddings=num_categories, embedding_dim=embed_dim)
        self.lin1 = torch.nn.Linear(256, 256)
        self.lin2 = torch.nn.Linear(256, 128)
        self.bn1 = torch.nn.BatchNorm1d(128)
        self.bn2 = torch.nn.BatchNorm1d(64)
        self.act1 = torch.nn.ReLU()
        self.act2 = torch.nn.ReLU()

    # Forward step of a model
    def forward(self, data):
        x, edge_index, batch = data.x, data.edge_index, data.batch

        item_id = x[:,:,0]
        category = x[:,:,1]


        emb_item = self.item_embedding(item_id).squeeze(1)
        emb_category = self.category_embedding(category).squeeze(1)
```

```python
        x = torch.cat([emb_item, emb_category], dim=1)
        # print(x.shape)
        x = F.relu(self.conv1(x, edge_index))
        # print(x.shape)
        r = self.pool1(x, edge_index, None, batch)
        # print(r)
        x, edge_index, _, batch, _, _ = self.pool1(x, edge_index, None, batch)
        x1 = torch.cat([gmp(x, batch), gap(x, batch)], dim=1)

        x = F.relu(self.conv2(x, edge_index))

        x, edge_index, _, batch, _, _ = self.pool2(x, edge_index, None, batch)
        x2 = torch.cat([gmp(x, batch), gap(x, batch)], dim=1)

        x = F.relu(self.conv3(x, edge_index))

        x, edge_index, _, batch, _, _ = self.pool3(x, edge_index, None, batch)
        x3 = torch.cat([gmp(x, batch), gap(x, batch)], dim=1)

        x = x1 + x2 + x3

        x = self.lin1(x)
        x = self.act1(x)
        x = self.lin2(x)
        x = F.dropout(x, p=0.5, training=self.training)
        x = self.act2(x)

        outputs = []
        for i in range(x.size(0)):
            output = torch.matmul(emb_item[data.batch == i], x[i,:])

            outputs.append(output)

        x = torch.cat(outputs, dim=0)
        x = torch.sigmoid(x)

        return x
```

## ▾ Обучение нейронной сверточной сети

```python
# Enable CUDA computing
device = torch.device('cuda')
model = Net().to(device)
# Choose optimizer and criterion for learning
optimizer = torch.optim.Adam(model.parameters(), lr=0.002)
crit = torch.nn.BCELoss()


# Train function
def train():
    model.train()

    loss_all = 0
    for data in train_loader:
        data = data.to(device)
        optimizer.zero_grad()
        output = model(data)

        label = data.y.to(device)
        loss = crit(output, label)
        loss.backward()
        loss_all += data.num_graphs * loss.item()
        optimizer.step()
    return loss_all / len(train_dataset)


# Evaluate result of a model
from sklearn.metrics import roc_auc_score
def evaluate(loader):
    model.eval()

    predictions = []
    labels = []

    with torch.no_grad():
```

```
        for data in loader:

            data = data.to(device)
            pred = model(data).detach().cpu().numpy()

            label = data.y.detach().cpu().numpy()
            predictions.append(pred)
            labels.append(label)

        predictions = np.hstack(predictions)
        labels = np.hstack(labels)

        return roc_auc_score(labels, predictions)
```

```
# Train a model                                          NUM_EPOCHS: 12
NUM_EPOCHS =  12 #@param { type: "integer" }
for epoch in tqdm(range(NUM_EPOCHS)):
    loss = train()
    train_acc = evaluate(train_loader)
    val_acc = evaluate(val_loader)
    test_acc = evaluate(test_loader)
    print('Epoch: {:03d}, Loss: {:.5f}, Train Auc: {:.5f}, Val Auc: {:.5f}, Test Auc: {:.5f}'.
          format(epoch, loss, train_acc, val_acc, test_acc))
```

```
  8%|█       | 1/12 [00:48<08:50, 48.23s/it]Epoch: 000, Loss: 0.66438, Train Auc: 0.53963, Val Auc: 0.54502, Test Au
 17%|█▋      | 2/12 [01:32<07:41, 46.10s/it]Epoch: 001, Loss: 0.49811, Train Auc: 0.58664, Val Auc: 0.57268, Test Au
 25%|██▌     | 3/12 [02:17<06:48, 45.42s/it]Epoch: 002, Loss: 0.44270, Train Auc: 0.61824, Val Auc: 0.58817, Test Aud
 33%|███▎    | 4/12 [03:02<06:00, 45.08s/it]Epoch: 003, Loss: 0.40889, Train Auc: 0.65620, Val Auc: 0.60315, Test Au
 42%|████▏   | 5/12 [03:46<05:14, 45.00s/it]Epoch: 004, Loss: 0.38101, Train Auc: 0.69724, Val Auc: 0.62240, Test Au
 50%|█████   | 6/12 [04:30<04:28, 44.70s/it]Epoch: 005, Loss: 0.35971, Train Auc: 0.71789, Val Auc: 0.62465, Test Au
 58%|█████▊  | 7/12 [05:15<03:43, 44.64s/it]Epoch: 006, Loss: 0.33764, Train Auc: 0.76569, Val Auc: 0.64367, Test Au
 67%|██████▋ | 8/12 [05:59<02:58, 44.53s/it]Epoch: 007, Loss: 0.32084, Train Auc: 0.79539, Val Auc: 0.64680, Test Au
 75%|███████▌| 9/12 [06:44<02:13, 44.47s/it]Epoch: 008, Loss: 0.30261, Train Auc: 0.82702, Val Auc: 0.65226, Test Au
 83%|████████▎| 10/12 [07:28<01:28, 44.37s/it]Epoch: 009, Loss: 0.28681, Train Auc: 0.85143, Val Auc: 0.65794, Test A
 92%|█████████▏| 11/12 [08:12<00:44, 44.46s/it]Epoch: 010, Loss: 0.26965, Train Auc: 0.87717, Val Auc: 0.65435, Test A
100%|██████████| 12/12 [08:56<00:00, 44.75s/it]Epoch: 011, Loss: 0.26058, Train Auc: 0.89207, Val Auc: 0.65470, Test Au
```

▼ Проверка результата с помощью примеров

```
# Подход №1 - из датасета
evaluate(DataLoader(test_dataset[25:45], batch_size=10))
```

```
    /usr/local/lib/python3.7/dist-packages/torch_geometric/deprecation.py:12: UserWarning: 'data.DataLoader' is deprecated,
      warnings.warn(out)
    0.8218390804597702
```

```
# Подход №2 - через создание сессии покупок
test_df = pd.DataFrame([
    [-1, 15219, 0],
    [-1, 15431, 0],
    [-1, 14371, 0],
    [-1, 15745, 0],
    [-2, 14594, 0],
    [-2, 16972, 11],
    [-2, 16943, 0],
    [-3, 17284, 0]
], columns=['session_id', 'item_id', 'category'])

test_data = transform_dataset(test_df, buy_item_dict)
test_data = DataLoader(test_data, batch_size=1)

with torch.no_grad():
    model.eval()
    for data in test_data:
        data = data.to(device)
        pred = model(data).detach().cpu().numpy()

        print(data, pred)
```

```
100%|██████████| 3/3 [00:00<00:00, 144.95it/s]DataBatch(x=[1, 1, 2], edge_index=[2, 0], y=[1], batch=[1], ptr=[2]) [0.8
DataBatch(x=[3, 1, 2], edge_index=[2, 2], y=[3], batch=[3], ptr=[2]) [0.05040454 0.01688893 0.01116091]
DataBatch(x=[4, 1, 2], edge_index=[2, 3], y=[4], batch=[4], ptr=[2]) [0.00638917 0.01556076 0.09219604 0.27488992]
```

```
/usr/local/lib/python3.7/dist-packages/torch_geometric/deprecation.py:12: UserWarning: 'data.DataLoader' is deprecated,
  warnings.warn(out)
```

Как видно из результатов, значение метрики *AUC = 72.5%*

В ходе работы были изменены следующие гиперпараметры: количество эпох (5->10), скорость обучение (0.001->0.002), количество сессий (50000->60000)

## **Вывод:**

При выполнении работы был подготовлен датасет графовых данных и подобрана модель и гиперпараметры обучения для получения качества AUC > 0.65. Удалось добиться значения метрики AUC = 0.725.