

딥러닝 모델 구현해 보기

학습 내용

- 정신 건강 설문 조사 데이터를 사용하여 개인의 우울증을 경험할 수 있는 요인을 탐색
- 첫번째 데이터 셋 : 자전거 공유 업체 시간대별 데이터
- 두번째 데이터 셋 : 타이타닉 데이터 셋(PyTorch)
- 세번째 데이터 셋 : 정신 건강 설문 조사 데이터

목차

01. 사전 환경 설치
02. 라이브러리 및 데이터 불러오기
03. 신경망 모델 정의
04. 예측 수행

01. 사전 환경 설치

목차로 이동하기

GPU 버전 PyTorch 설치

```
# 가상환경 만들기
conda create --name gpuDL python=3.10

# Jupyter Notebook
conda install -c conda-forge notebook jupyter

# 추가 설치
pip install pandas scikit-learn

# PyTorch 설치하기
# 01 Anaconda 사용 시
conda install pytorch=2.5.1 torchvision torchaudio pytorch-
cuda=12.1 -c pytorch -c nvidia

# 02 pip 사용 시
pip3 install torch==2.5.1 torchvision torchaudio --index-url
https://download.pytorch.org/whl/cu121
```

```
In [4]: import torch
import sys
import numpy
import torch
```

```

print("Python version:", sys.version)
print("NumPy version:", numpy.__version__)
print("PyTorch version:", torch.__version__)

# CUDA 사용 가능 여부 확인
print(torch.cuda.is_available())

# 사용 가능한 GPU 장치 수 확인
print(torch.cuda.device_count())

```

Python version: 3.11.10 | packaged by Anaconda, Inc. | (main, Oct 3 2024, 07:22:26) [MSC v.1929 64 bit (AMD64)]
 NumPy version: 2.1.3
 PyTorch version: 2.5.1+cpu
 False
 0

PyTorch 버전	권장 CUDA 버전
1.13.x	11.7
2.0.x	11.8
2.1.x --	12.1

```

In [8]: import torch
import sklearn

print(torch.__version__)
print(sklearn.__version__)

```

2.5.1+cpu
 1.5.2

02. 라이브러리 및 데이터 불러오기

[목차로 이동하기](#)

```

In [7]: import numpy as np
import pandas as pd
import torch
import torch.nn as nn
import torch.optim as optim
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.impute import SimpleImputer

```

```

In [12]: import os
os.getcwd()

```

Out[12]: 'd:\\github\\DeepLearning_Basic_Class'

```

In [44]: # 시드 고정
torch.manual_seed(42)
np.random.seed(42)

# 1. 데이터 준비
# 상대 경로로 데이터 로드

```

```
train = pd.read_csv('./datasets/health_mental_24/train.csv')
test = pd.read_csv('./datasets/health_mental_24/test.csv')
sub = pd.read_csv('./datasets/health_mental_24/sample_submission.csv')

# 데이터 shape 확인
print("훈련 데이터 shape:", train.shape)
print("테스트 데이터 shape:", test.shape)

# 데이터 정보 확인
print("\n훈련 데이터 정보:")
print(train.info())

print("\n테스트 데이터 정보:")
print(test.info())
```

훈련 데이터 shape: (140700, 20)
테스트 데이터 shape: (93800, 19)

훈련 데이터 정보:

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 140700 entries, 0 to 140699  
Data columns (total 20 columns):
```

#	Column	Non-Null Count	Dtype
0	id	140700 non-null	int64
1	Name	140700 non-null	object
2	Gender	140700 non-null	object
3	Age	140700 non-null	float64
4	City	140700 non-null	object
5	Working Professional or Student	140700 non-null	object
6	Profession	104070 non-null	object
7	Academic Pressure	27897 non-null	float64
8	Work Pressure	112782 non-null	float64
9	CGPA	27898 non-null	float64
10	Study Satisfaction	27897 non-null	float64
11	Job Satisfaction	112790 non-null	float64
12	Sleep Duration	140700 non-null	object
13	Dietary Habits	140696 non-null	object
14	Degree	140698 non-null	object
15	Have you ever had suicidal thoughts ?	140700 non-null	object
16	Work/Study Hours	140700 non-null	float64
17	Financial Stress	140696 non-null	float64
18	Family History of Mental Illness	140700 non-null	object
19	Depression	140700 non-null	int64

dtypes: float64(8), int64(2), object(10)

memory usage: 21.5+ MB

None

테스트 데이터 정보:

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 93800 entries, 0 to 93799  
Data columns (total 19 columns):
```

#	Column	Non-Null Count	Dtype
0	id	93800 non-null	int64
1	Name	93800 non-null	object
2	Gender	93800 non-null	object
3	Age	93800 non-null	float64
4	City	93800 non-null	object
5	Working Professional or Student	93800 non-null	object
6	Profession	69168 non-null	object
7	Academic Pressure	18767 non-null	float64
8	Work Pressure	75022 non-null	float64
9	CGPA	18766 non-null	float64
10	Study Satisfaction	18767 non-null	float64
11	Job Satisfaction	75026 non-null	float64
12	Sleep Duration	93800 non-null	object
13	Dietary Habits	93795 non-null	object
14	Degree	93798 non-null	object
15	Have you ever had suicidal thoughts ?	93800 non-null	object
16	Work/Study Hours	93800 non-null	float64
17	Financial Stress	93800 non-null	float64
18	Family History of Mental Illness	93800 non-null	object

dtypes: float64(8), int64(1), object(10)

memory usage: 13.6+ MB
None

In [45]: `train.head()`

Out[45]:

	id	Name	Gender	Age	City	Working Professional or Student	Profession	Academic Pressure	Pre
0	0	Aaradhya	Female	49.0	Ludhiana	Working Professional	Chef	NaN	
1	1	Vivan	Male	26.0	Varanasi	Working Professional	Teacher	NaN	
2	2	Yuvraj	Male	33.0	Visakhapatnam	Student	NaN	5.0	
3	3	Yuvraj	Male	22.0	Mumbai	Working Professional	Teacher	NaN	
4	4	Rhea	Female	30.0	Kanpur	Working Professional	Business Analyst	NaN	

In [46]: `test.head(3)`

Out[46]:

	id	Name	Gender	Age	City	Working Professional or Student	Profession	Academic Pressure	
0	140700	Shivam	Male	53.0	Visakhapatnam	Working Professional	Judge	NaN	
1	140701	Sanya	Female	58.0	Kolkata	Working Professional	Educational Consultant	NaN	
2	140702	Yash	Male	53.0	Jaipur	Working Professional	Teacher	NaN	

In [47]: `sub.head()`

```
Out[47]:
```

	id	Depression
0	140700	0
1	140701	0
2	140702	0
3	140703	0
4	140704	0

```
In [48]: train.columns
```

```
Out[48]: Index(['id', 'Name', 'Gender', 'Age', 'City',
               'Working Professional or Student', 'Profession', 'Academic Pressure',
               'Work Pressure', 'CGPA', 'Study Satisfaction', 'Job Satisfaction',
               'Sleep Duration', 'Dietary Habits', 'Degree',
               'Have you ever had suicidal thoughts ?', 'Work/Study Hours',
               'Financial Stress', 'Family History of Mental Illness', 'Depression'],
              dtype='object')
```

```
In [49]: # 필요한 피쳐 선택
features = ['Age', 'Work/Study Hours', 'Financial Stress', 'Gender']
target = 'Depression'

# 성별 인코딩
train['Gender'] = train['Gender'].map({'Male': 0, 'Female': 1})

# 결측값 처리
imputer = SimpleImputer(strategy='median')
X = imputer.fit_transform(train[features])
y = train[target].values
```

```
In [50]: # 스케일링
scaler = StandardScaler()
X = scaler.fit_transform(X)

# 데이터 분할
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)

# NumPy to PyTorch Tensor 변환
# X_train이라는 NumPy 배열을 PyTorch의 FloatTensor로 변환.
# FloatTensor는 32비트 부동 소수점 숫자로 구성된 텐서를 생성.
# 이 변환은 PyTorch 모델에서 데이터를 처리할 수 있도록 준비하는 단계
X_train = torch.FloatTensor(X_train)
X_test = torch.FloatTensor(X_test)
y_train = torch.FloatTensor(y_train).unsqueeze(1)
y_test = torch.FloatTensor(y_test).unsqueeze(1)
```

03. 신경망 모델 정의

목차로 이동하기

- 딥러닝의 이해를 위해 일부 특징(변수)만 지정하였음.
- 이미지를 사용할 때는 지정된 이미지 전체를 입력 데이터로 사용하는 경우가 대부분.

```
In [53]: # 2. 신경망 모델 정의
model = nn.Sequential(
    nn.Linear(4, 8),
    nn.ReLU(),
    nn.Dropout(0.3),
    nn.Linear(8, 4),
    nn.ReLU(),
    nn.Dropout(0.3),
    nn.Linear(4, 1),
    nn.Sigmoid()
)

# 3. 모델 학습
# 손실 함수와 옵티마이저 정의
criterion = nn.BCELoss() # 이진 분류 손실 함수
optimizer = optim.Adam(model.parameters(), lr=0.001)
```

```
In [54]: # 학습 진행
epochs = 200
for epoch in range(epochs):
    # 순전파
    outputs = model(X_train)
    loss = criterion(outputs, y_train)

    # 정확도 계산
    with torch.no_grad():
        # 이진 분류의 경우
        predicted = (outputs > 0.5).float()
        accuracy = (predicted == y_train).float().mean()

    # 역전파
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

    # 20번마다 손실과 정확도 출력
    if (epoch + 1) % 20 == 0:
        print(f'Epoch [{epoch+1}/{epochs}], '
              f'Loss: {loss.item():.4f}, '
              f'Accuracy: {accuracy.item():.4f}')
```

```
Epoch [20/200], Loss: 0.5721, Accuracy: 0.8186
Epoch [40/200], Loss: 0.5589, Accuracy: 0.8186
Epoch [60/200], Loss: 0.5433, Accuracy: 0.8186
Epoch [80/200], Loss: 0.5241, Accuracy: 0.8186
Epoch [100/200], Loss: 0.5014, Accuracy: 0.8186
Epoch [120/200], Loss: 0.4769, Accuracy: 0.8186
Epoch [140/200], Loss: 0.4549, Accuracy: 0.8186
Epoch [160/200], Loss: 0.4341, Accuracy: 0.8186
Epoch [180/200], Loss: 0.4145, Accuracy: 0.8186
Epoch [200/200], Loss: 0.3987, Accuracy: 0.8186
```

```
In [55]: # 4. 모델 평가
model.eval() # 평가 모드
with torch.no_grad():
    test_outputs = model(X_test)
    predicted = (test_outputs > 0.5).float()
    accuracy = (predicted == y_test).float().mean()
    print(f'Test Accuracy: {accuracy.item():.4f}')
```

Test Accuracy: 0.8168

04. 새로운 데이터로 예측

목차로 이동하기

```
In [56]: # 필요한 피처 선택
features = ['Age', 'Work/Study Hours', 'Financial Stress', 'Gender']

# Gender 인코딩
test['Gender'] = test['Gender'].map({'Male': 0, 'Female': 1})

# 선택한 피처로 test 데이터 준비
X_predict = test[features]

# 결측값 처리
X_predict = scaler.transform(X_predict)

# 스케일링
X_predict = scaler.transform(X_predict)

# PyTorch Tensor로 변환
X_predict = torch.FloatTensor(X_predict)

# 예측 수행
model.eval()
with torch.no_grad():
    pred = model(X_predict)
    pred = (pred > 0.5).float()

# submission 파일에 예측값 저장
sub['Depression'] = pred.numpy()

sub.to_csv('./datasets/health_mental_24/sub01.csv', index=False)
print("예측 완료 및 sub 파일 저장.")
```

c:\Users\daniel_wj\anaconda3\envs\gpuDL\Lib\site-packages\sklearn\base.py:486: UserWarning: X has feature names, but StandardScaler was fitted without feature names

warnings.warn(
예측 완료 및 sub 파일 저장.

추가 실습

- 여러개의 특징을 선택 및 신경망의 뉴런 추가 등으로 성능을 개선시켜 보자.