

# 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

```
In [1]: 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
```

```
In [2]: 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 [3]: 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 [4]: import os
os.getcwd()
```

```
Out[4]: 'd:\\github\\DeepLearning_Basic_Class'
```

```
In [5]: # 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 [6]: `train.columns`

Out[6]: 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 [7]: import numpy as np
import pandas as pd
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import LabelEncoder, StandardScaler

# 필요한 피처 선택
features = [
    'Age', 'Work/Study Hours', 'Financial Stress', # 숫자형 변수
    'Gender', 'Working Professional or Student', 'City',
    'Sleep Duration', 'Dietary Habits', 'Degree', # 범주형 변수
    'Academic Pressure', 'Work Pressure', 'CGPA',
    'Study Satisfaction', 'Job Satisfaction' # 결측치 있는 숫자형 변수
]
target = 'Depression'
```

```
In [8]: def create_feature_engineering(df):
# 데이터프레임 복사 (원본 데이터 보존)
df = df.copy()

# 1. Stress Index 생성 (결측값 처리 포함)
df['Stress_Index'] = (
    df['Academic Pressure'].fillna(df['Academic Pressure'].median()) +
    df['Work Pressure'].fillna(df['Work Pressure'].median())
) / 2

# 2. 연령대 그룹화 (결측값 처리)
df['Age_Group'] = pd.cut(
    df['Age'].fillna(df['Age'].median()),
    bins=[0, 20, 30, 40, 50, 100],
    labels=['Teen', 'Young Adult', 'Adult', 'Middle Age', 'Senior']
)

# 3. 수면 시간과 스트레스 상관관계 (범주형 변수 처리)
sleep_mapping = {
    'Very Short': 0,
    'Short': 1,
    'Normal': 2,
    'Long': 3,
    'Very Long': 4
}
df['Sleep_Duration_Numeric'] = df['Sleep Duration'].map(sleep_mapping).fillna(0)
df['Sleep_Stress_Interaction'] = df['Sleep_Duration_Numeric'] * df['Stress_Index']

# 4. 학업/직업 시간 비율 (0으로 나누기 방지)
df['Work_Study_Ratio'] = df['Work/Study Hours'] / (df['Age'] + 1) # 0으로 나눌 수 없게 하기

# 5. 추가 파생 변수들
```

```

# 만족도 지표
df['Satisfaction_Index'] = (
    df['Study Satisfaction'].fillna(df['Study Satisfaction'].median()) +
    df['Job Satisfaction'].fillna(df['Job Satisfaction'].median())
) / 2

# 6. 이상값 처리
def handle_outliers(series):
    Q1 = series.quantile(0.25)
    Q3 = series.quantile(0.75)
    IQR = Q3 - Q1
    upper_bound = Q3 + 1.5 * IQR
    lower_bound = Q1 - 1.5 * IQR
    return series.clip(lower_bound, upper_bound)

# 이상값 처리할 변수들
outlier_columns = ['Stress_Index', 'Work_Study_Ratio', 'Satisfaction_Index']
for col in outlier_columns:
    df[col] = handle_outliers(df[col])

return df

```

```

In [12]: # 전처리 함수 정의
def preprocess_data(df, gubun='train'):
    # 데이터프레임 복사
    data = df.copy()

    # 특징 엔지니어링 적용
    data = create_feature_engineering(data)

    # 필요한 피처 선택 (특징 엔지니어링 후 추가된 피처 포함)
    features = [
        'Age', 'Work/Study Hours', 'Financial Stress', # 기존 숫자형 변수
        'Stress_Index', 'Work_Study_Ratio', 'Satisfaction_Index', # 새로 추가된
        'Gender', 'Working Professional or Student', 'City',
        'Sleep Duration', 'Dietary Habits', 'Degree', # 범주형 변수
        'Academic Pressure', 'Work Pressure', 'CGPA',
        'Study Satisfaction', 'Job Satisfaction', # 기존 숫자형 변수
        'Sleep_Duration_Numeric', 'Sleep_Stress_Interaction' # 새로 추가된 상호
    ]

    # 범주형 변수 인코딩
    categorical_features = [
        'Gender',
        'Working Professional or Student',
        'City',
        'Sleep Duration',
        'Dietary Habits',
        'Degree'
        # 'Age_Group' # 새로 추가된 범주형 변수
    ]

    # 라벨 인코더 초기화
    le = LabelEncoder()

    # 범주형 변수 인코딩
    for col in categorical_features:
        # NaN 값을 문자열로 변환 후 인코딩
        data[col] = data[col].fillna('Unknown')

```

```

data[col] = le.fit_transform(data[col].astype(str))

# 숫자형 변수 분리
numeric_features = [
    'Age', 'Work/Study Hours', 'Financial Stress',
    'Academic Pressure', 'Work Pressure', 'CGPA',
    'Study Satisfaction', 'Job Satisfaction',
    'Stress_Index', 'Work_Study_Ratio', 'Satisfaction_Index',
    'Sleep_Duration_Numeric', 'Sleep_Stress_Interaction'
]

# 결측값 처리
imputer = SimpleImputer(strategy='median')

# 숫자형 변수 결측값 대체
numeric_data = imputer.fit_transform(data[numeric_features])

# 스케일러로 숫자형 변수 표준화
scaler = StandardScaler()
scaled_numeric_data = scaler.fit_transform(numeric_data)

# 범주형 변수와 숫자형 변수 결합
X = np.column_stack([
    scaled_numeric_data, # 표준화된 숫자형 변수
    data[categorical_features].values # 인코딩된 범주형 변수
])

if gubun=="train":
    # 타겟 변수 추출
    y = data[target].values
    return X, y
else:
    return X

# 전처리 적용
X, y = preprocess_data(train)
print(X.shape)

```

(140700, 19)

```

In [13]: # 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 [17]: # 2. 신경망 모델 정의
model = nn.Sequential(
    nn.Linear(19, 32),
    nn.ReLU(),
    nn.Dropout(0.3),
    nn.Linear(32, 16),
    nn.ReLU(),
    nn.Dropout(0.3),
    nn.Linear(16, 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 [18]: # 학습 진행
epochs = 500
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 [20/500], Loss: 0.5907, Accuracy: 0.8365
Epoch [40/500], Loss: 0.4996, Accuracy: 0.8443
Epoch [60/500], Loss: 0.3933, Accuracy: 0.8721
Epoch [80/500], Loss: 0.3091, Accuracy: 0.8960
Epoch [100/500], Loss: 0.2612, Accuracy: 0.9075
Epoch [120/500], Loss: 0.2374, Accuracy: 0.9104
Epoch [140/500], Loss: 0.2263, Accuracy: 0.9126
Epoch [160/500], Loss: 0.2215, Accuracy: 0.9140
Epoch [180/500], Loss: 0.2174, Accuracy: 0.9146
Epoch [200/500], Loss: 0.2147, Accuracy: 0.9147
Epoch [220/500], Loss: 0.2125, Accuracy: 0.9155
Epoch [240/500], Loss: 0.2117, Accuracy: 0.9157
Epoch [260/500], Loss: 0.2101, Accuracy: 0.9169
Epoch [280/500], Loss: 0.2088, Accuracy: 0.9179
Epoch [300/500], Loss: 0.2077, Accuracy: 0.9177
Epoch [320/500], Loss: 0.2075, Accuracy: 0.9176
Epoch [340/500], Loss: 0.2059, Accuracy: 0.9179
Epoch [360/500], Loss: 0.2052, Accuracy: 0.9184
Epoch [380/500], Loss: 0.2049, Accuracy: 0.9186
Epoch [400/500], Loss: 0.2039, Accuracy: 0.9184
Epoch [420/500], Loss: 0.2041, Accuracy: 0.9188
Epoch [440/500], Loss: 0.2039, Accuracy: 0.9181
Epoch [460/500], Loss: 0.2038, Accuracy: 0.9184
Epoch [480/500], Loss: 0.2028, Accuracy: 0.9192
Epoch [500/500], Loss: 0.2023, Accuracy: 0.9190
```

```
In [19]: # 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: 0.9221

## 04. 새로운 데이터로 예측

목차로 이동하기

```
In [20]: # 필요한 피처 선택
# features = [
#     'Age', 'Work/Study Hours', 'Financial Stress', # 기존 숫자형 변수
#     'Stress_Index', 'Work_Study_Ratio', 'Satisfaction_Index', # 새로 추가
#     'Gender', 'Working Professional or Student', 'City',
#     'Sleep Duration', 'Dietary Habits', 'Degree', # 범주형 변수
#     'Academic Pressure', 'Work Pressure', 'CGPA',
#     'Study Satisfaction', 'Job Satisfaction', # 기존 숫자형 변수
#     'Sleep_Duration_Numeric', 'Sleep_Stress_Interaction' # 새로 추가된 상
# ]

# 선택한 피처로 test 데이터 준비
X_predict = preprocess_data(test, gubun='test')

# 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/sub03.csv', index=False)
print("예측 완료 및 sub 파일 저장.")

```

예측 완료 및 sub 파일 저장.

## 앙상블 모델

```

In [22]: # 데이터 분할
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)

from sklearn.ensemble import RandomForestClassifier

model = RandomForestClassifier()
model.fit(X_train, y_train)
model.score(X_test, y_test)

```

Out[22]: 0.9212508884150675

```

In [23]: # 선택한 피처로 test 데이터 준비
X_predict = preprocess_data(test, gubun='test')

# 예측 수행
pred = model.predict(X_predict)

# submission 파일에 예측값 저장
sub['Depression'] = pred

# submission 파일 저장
sub.to_csv('./datasets/health_mental_24/sub_rf.csv', index=False)
print("예측 완료 및 submission 파일 저장됨")

```

예측 완료 및 submission 파일 저장됨

## 참고 NOTEBOOK

- <https://www.kaggle.com/code/igorvolianiuk/mental-health-catboost-first-place>