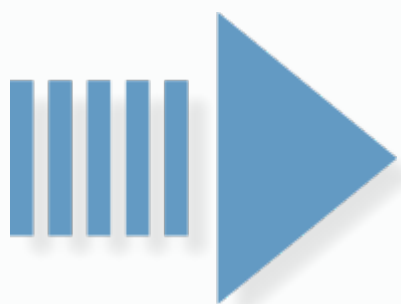


Shape 확인

```
1 import pandas as pd
2 import numpy as np
3 path = '/content/drive/MyDrive/airline_dataset/'
4 train = pd.read_csv(path+"train.csv")
5 test = pd.read_csv(path+'test.csv')
6 sample = (path+'sample_submission.csv')
7
8
9 print(f'train set은 {train.shape[1]} 개의 feature를 가진 {train.shape[0]} 개의 데이터 샘플로 이루어져 있습니다.')
10 print(f'test set은 {test.shape[1]} 개의 feature를 가진 {test.shape[0]} 개의 데이터 샘플로 이루어져 있습니다.')
11 test.head()
```



train set은 24 개의 feature를 가진 3000 개의 데이터 샘플로 이루어져 있습니다.
test set은 23 개의 feature를 가진 2000 개의 데이터 샘플로 이루어져 있습니다.

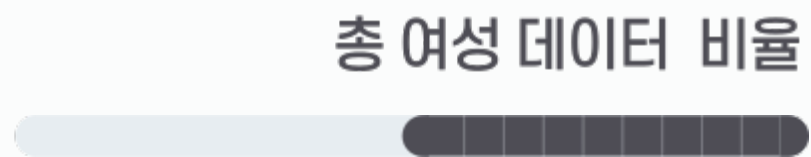
※ train set은 target을 포함하므로 24의 feature를 가집니다.

결측치 확인

```
▶ 1 def check_missing_col(dataframe):
2     missing_col = []
3     for col in dataframe.columns:
4         missing_values = sum(dataframe[col].isna())
5         is_missing = True if missing_values >= 1 else False
6         if is_missing:
7             print(f'결측치가 있는 컬럼은: {col} 입니다')
8             print(f'해당 컬럼에 총 {missing_values} 개의 결측치가 존재합니다.')
9             missing_col.append([col, dataframe[col].dtype])
10    if missing_col == []:
11        print('결측치가 존재하지 않습니다')
12    return missing_col
13
14 missing_col = check_missing_col(train)
15 missing_col = check_missing_col(test)
16
17 train.head()
18
```



결측치가 존재하지 않습니다
결측치가 존재하지 않습니다



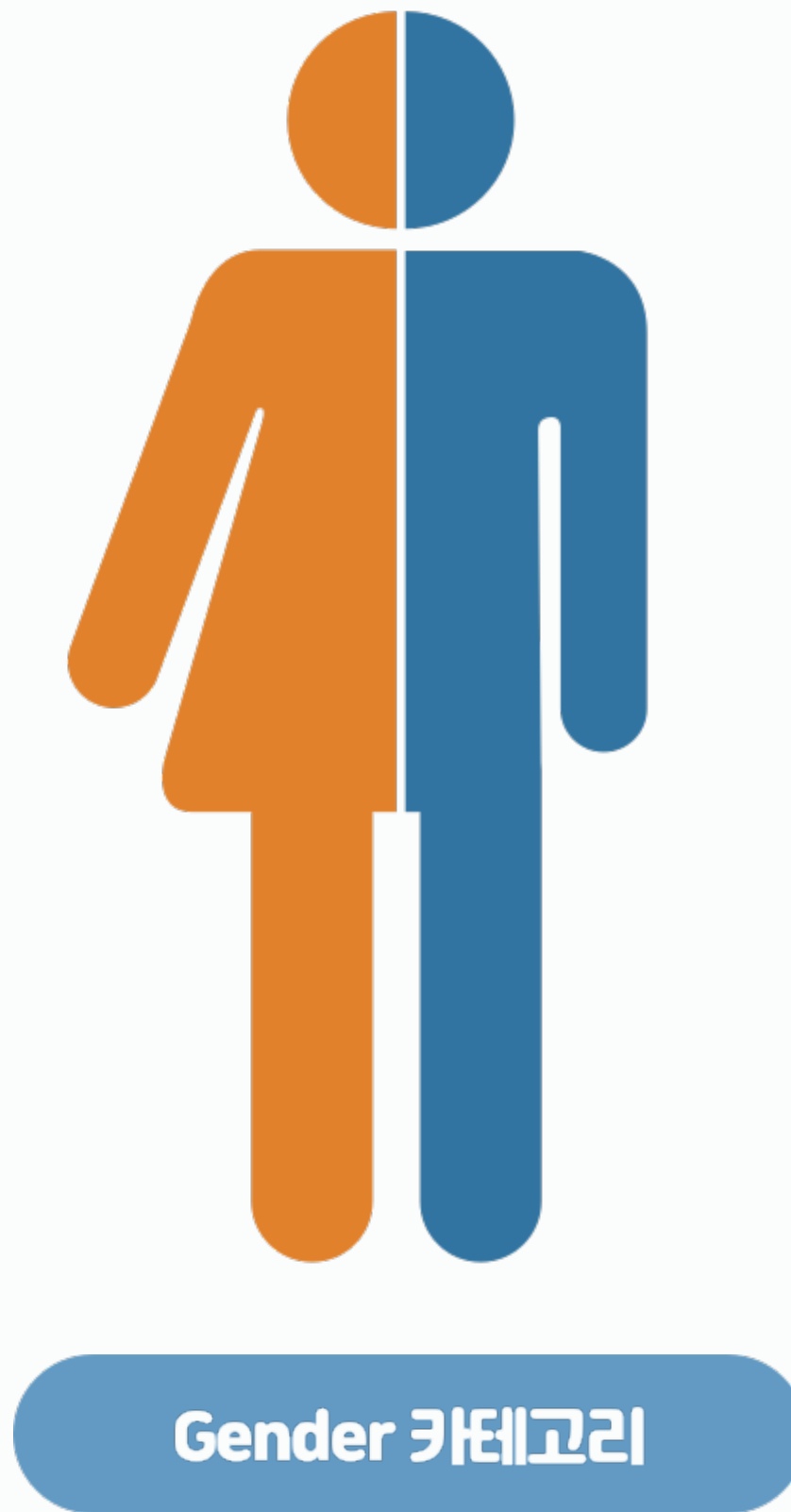
50.1%



60.4%



36.9%



총 남성 데이터 비율



49.9%



39.6%



63.1%

데이터 전처리

- 평균 지연 시간 컬럼을 생성

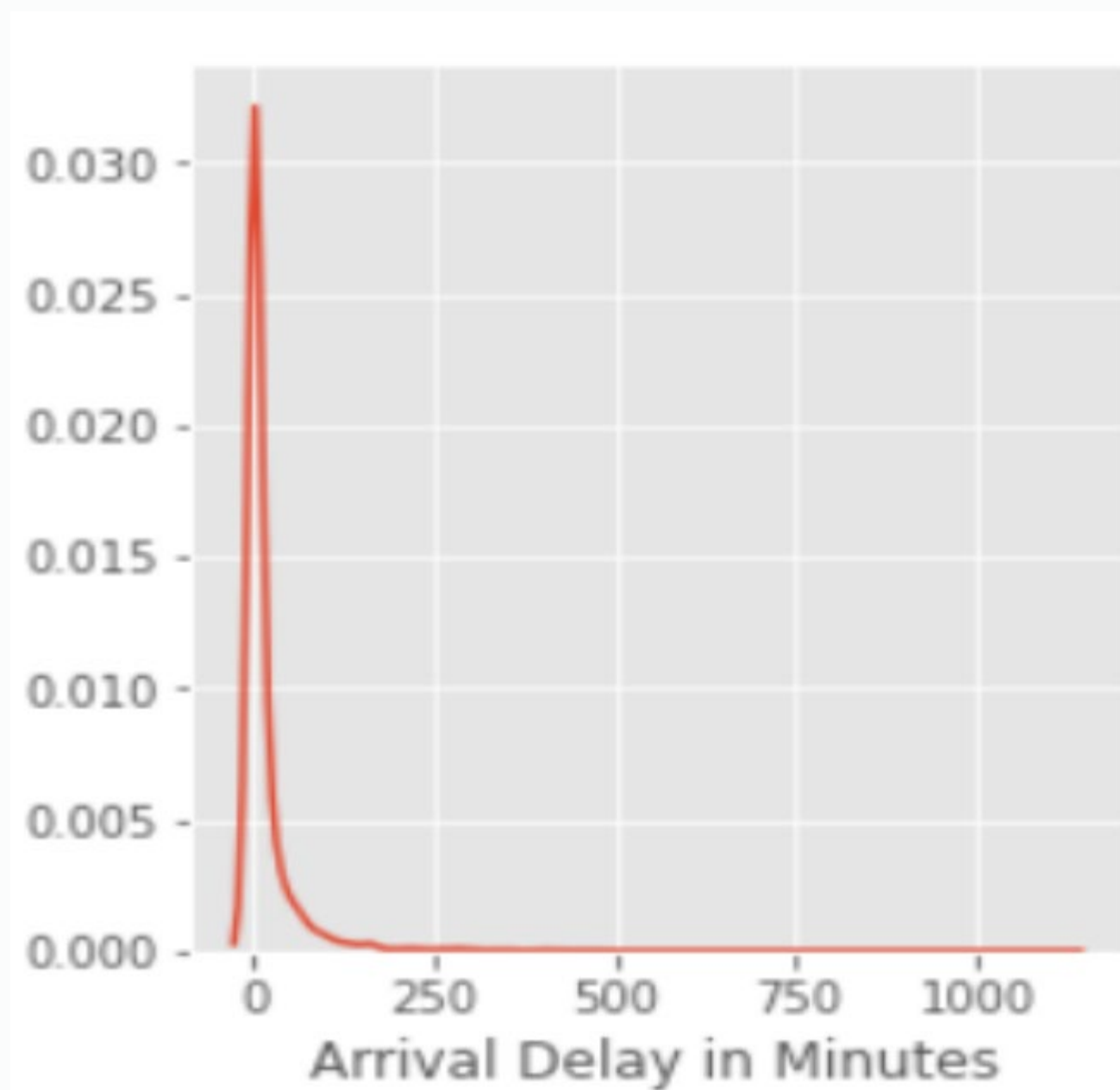
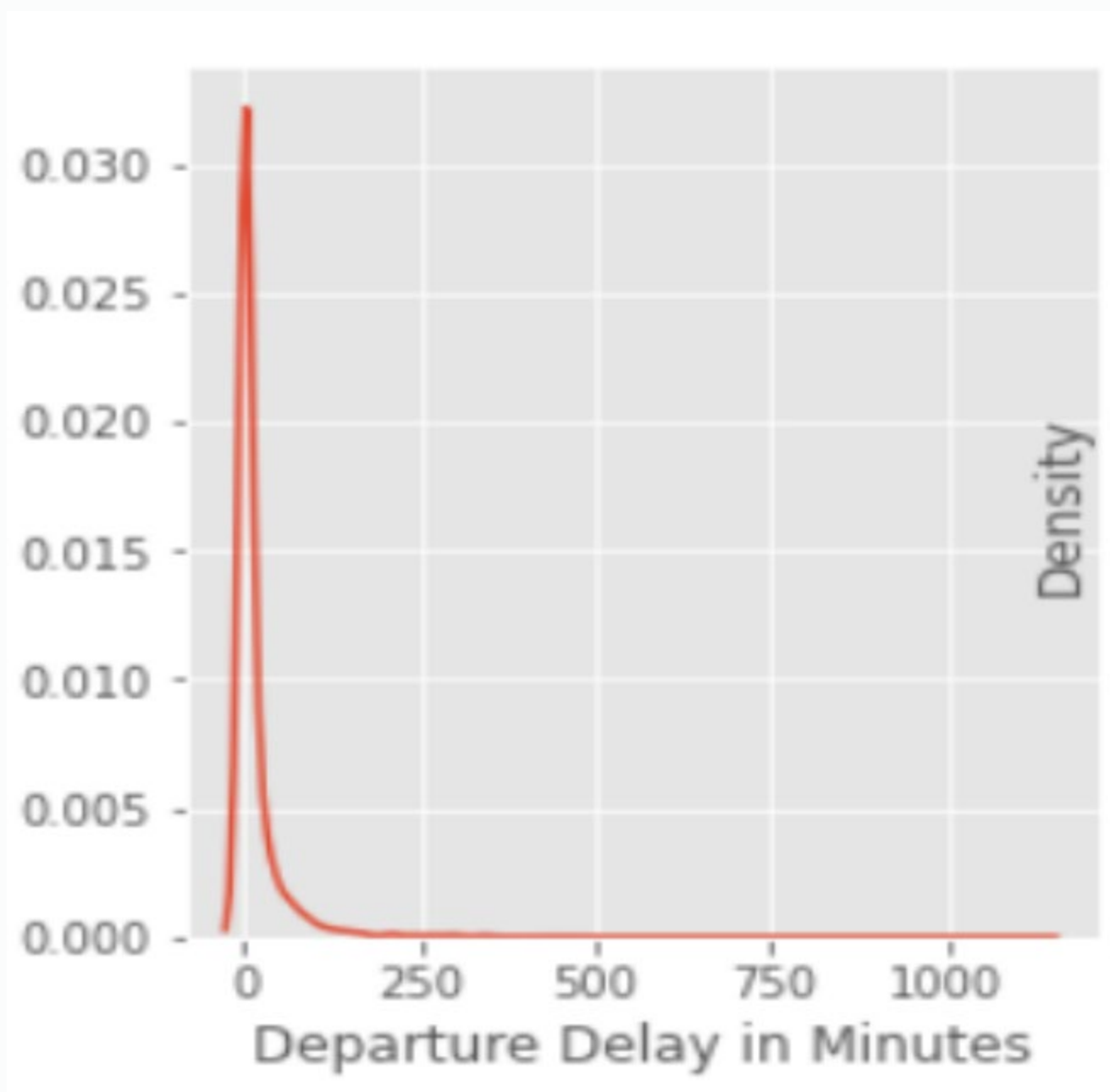
```
34 train["Mean Delay in Minutes"] = (train["Departure Delay in Minutes"] + train['Arrival Delay in Minutes']) / 2
35 test["Mean Delay in Minutes"] = (test["Departure Delay in Minutes"] + test['Arrival Delay in Minutes']) / 2
36
37
```

- 불필요한 컬럼은 제거

```
47 train_x = train.drop(["id", "Departure Delay in Minutes", "Arrival Delay in Minutes", "Food and drink", 'target'], axis=1)
48 train_y = train.target
49
50 test_x = test.drop(["id", "Departure Delay in Minutes", "Arrival Delay in Minutes", "Food and drink"], axis=1)
51
52
```

데이터 전처리-로그(log) 변환

- 큰 값은 엄청 큰 우측 꼬리가 긴 분포형태
→ log 변환을 이용하여 치우쳐진 정도를 줄임

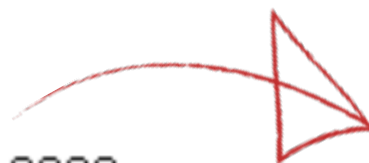


데이터 전처리-로그(log) 변환

- 0은 로그변환이 안되므로 전체 값에 1을 더한 뒤 로그변환 하는 `log1p` 함수를 사용해야 한다

```
1 #log 변환 실행
2 train['Mean Delay in Minutes'] = np.log1p(train['Mean Delay in Minutes'])
3
4 # test 데이터에도 변환 실행
5 test['Mean Delay in Minutes'] = np.log1p(test['Mean Delay in Minutes'])
6
```

Departure Delay in Minutes Skew : 9.190139679910239
Arrival Delay in Minutes Skew : 8.887761727831762



Departure Delay in Minutes Skew : 0.9302111175258293
Arrival Delay in Minutes Skew : 0.8979015577156512

↳ 모든 skew(왜도) 가 1 이하로 내려간 모습

라벨 인코딩

라벨인코딩을 하기 위한 dictionary map 생성 함수

```
1 def make_label_map(dataframe):
2     label_maps = {}
3     for col in dataframe.columns:
4         if dataframe[col].dtype=='object':
5             label_map = {'unknown':0}
6             for i, key in enumerate(dataframe[col].unique()):
7                 label_map[key] = i+1 ← 새로 등장하는 유니크 값들에 대해 1부터 1씩 증가시켜 키값을 부여해줍니다.
8             label_maps[col] = label_map
9     return label_maps
```

train 데이터 라벨 인코딩

각 범주형 변수에 인코딩 값을 부여하는 함수

```
12 def label_encoder(dataframe, label_map):
13     for col in dataframe.columns:
14         if dataframe[col].dtype=='object':
15             dataframe[col] = dataframe[col].map(label_map[col])
16             dataframe[col] = dataframe[col].fillna(label_map[col]['unknown'])
17     return dataframe
```

```
20 label_map = make_label_map(train) # train 사용해 label map 생성
21 train_x = label_encoder(train, label_map) # train 라벨 인코딩
22
23
24 train_x.head()
25
26
```


라벨 인코딩

| | id | Gender | Customer Type | Age | Type of Travel | Class | Flight Distance | Seat comfort |
|---|----|--------|-------------------|-----|-----------------|----------|-----------------|--------------|
| 0 | 1 | Female | disloyal Customer | 22 | Business travel | Eco | 1599 | 3 |
| 1 | 2 | Female | Loyal Customer | 37 | Business travel | Business | 2810 | 2 |
| 2 | 3 | Male | Loyal Customer | 46 | Business travel | Business | 2622 | 1 |
| 3 | 4 | Female | disloyal Customer | 24 | Business travel | Eco | 2348 | 3 |
| 4 | 5 | Female | Loyal Customer | 58 | Business travel | Business | 105 | 3 |

- 라벨 인코딩 전

문자열 형식 → 데이터 숫자로 표현

| | Gender | Customer Type | Age | Type of Travel | Class | Flight Distance | Seat comfort | De |
|---|--------|---------------|----------|----------------|-------|-----------------|--------------|----|
| 0 | 1 | 1 | 0.205479 | 1 | 1 | 0.226501 | 3 | |
| 1 | 1 | 2 | 0.410959 | 1 | 2 | 0.403807 | 2 | |
| 2 | 2 | 2 | 0.534247 | 1 | 2 | 0.376281 | 1 | |
| 3 | 1 | 1 | 0.232877 | 1 | 1 | 0.336164 | 3 | |
| 4 | 1 | 2 | 0.698630 | 1 | 2 | 0.007760 | 3 | |

- 라벨 인코딩 후

min-max 정규화

- 수치형 데이터의 값을 0~1 사이의 값으로 변환
-데이터 값의 크기를 줄이고 분산을 줄여 모델이 데이터를 이상하게 해석하는 것을 방지

```
1 from sklearn.preprocessing import MinMaxScaler
2
3 num_features = ['Age', 'Flight Distance', 'Mean Delay in Minutes']
4
5 scaler = MinMaxScaler()
6 train_x[num_features] = scaler.fit_transform(train_x[num_features])
7 train_x.head()
```



| | id | Gender | Customer Type | Age | Type of Travel | Class | Flight Distance | Seat comfort | Departure/Arrival time convenient | Food and drink | ... | On-board service | Leg room service | Baggage handling | Checkin service | Cleanliness | Onlin boarding |
|---|----|--------|---------------|----------|----------------|-------|-----------------|--------------|-----------------------------------|----------------|-----|------------------|------------------|------------------|-----------------|-------------|----------------|
| 0 | 1 | 1 | 1 | 0.205479 | 1 | 1 | 0.226501 | 3 | 0 | 3 | ... | 5 | 4 | 4 | 4 | 5 | |
| 1 | 2 | 1 | 2 | 0.410959 | 1 | 2 | 0.403807 | 2 | 4 | 4 | ... | 5 | 4 | 2 | 1 | 5 | |
| 2 | 3 | 2 | 2 | 0.534247 | 1 | 2 | 0.376281 | 1 | 1 | 1 | ... | 4 | 4 | 4 | 5 | 4 | |
| 3 | 4 | 1 | 1 | 0.232877 | 1 | 1 | 0.336164 | 3 | 3 | 3 | ... | 2 | 4 | 5 | 3 | 4 | |
| 4 | 5 | 1 | 2 | 0.698630 | 1 | 2 | 0.007760 | 3 | 3 | 3 | ... | 4 | 4 | 4 | 4 | 4 | |

5 rows x 25 columns

Batch_size , epoch

- batch_size = 한 번 연산할 때 들어가는 데이터의 크기
- epoch = 학습 횟수

```
1 #Define training hyperparameters.  
2 batch_size = 50  
3 num_epochs = 1000  
4  
5 #Calculate some other hyperparameters based on data.  
6 batch_no = len(train_x) // batch_size #batches  
7 cols=train_x.shape[1] #Number of columns in input matrix  
8 n_output=1  
9
```

Net 설계

```
1 #####
2 import torch
3 import torch.nn as nn
4 #Create the model
5 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
6 # Assume that we are on a CUDA machine, then this should print a CUDA device:
7 print("Executing the model on :",device)
8 class Net(torch.nn.Module):
9     def __init__(self,cols, n_output):
10         super().__init__()
11         self.linear1=nn.Linear(cols,100)
12         self.linear2=nn.Linear(100,50)
13         self.linear3=nn.Linear(50,30)
14         self.linear4=nn.Linear(30,10)
15         self.linear5=nn.Linear(10, n_output)
16
17         #self.dropout = nn.Dropout(p=0.2)
18         self.batchnorm1 = nn.BatchNorm1d(100)
19         self.batchnorm2 = nn.BatchNorm1d(50)
20         self.batchnorm3 = nn.BatchNorm1d(30)
21         self.batchnorm4 = nn.BatchNorm1d(10)
22
23         self.relu=nn.ReLU()
24
25         self.sigmoid=nn.Sigmoid()
26
```

```
26
27     def forward(self,X):
28         out=self.linear1(X)
29         out=self.batchnorm1(out)
30         out=self.relu(out)
31
32         out=self.linear2(out)
33         out=self.batchnorm2(out)
34         out=self.relu(out)
35
36         out=self.linear3(out)
37         out=self.batchnorm3(out)
38         out=self.relu(out)
39
40         out=self.linear4(out)
41         out=self.batchnorm4(out)
42         out=self.relu(out)
43
44
45         out=self.linear5(out)
46         #out=self.dropout(out)
47         out=self.sigmoid(out)
48         return out
49
50 model = Net(cols, n_output)
```



```

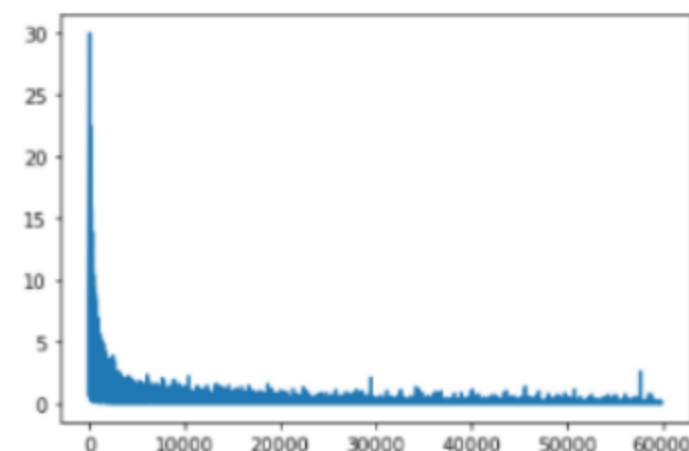
1 import torch.nn.functional as F
2 from sklearn.utils import shuffle
3 from torch.autograd import Variable
4 from matplotlib import pyplot as plt
5 running_loss = 0.0
6 losses=[]
7 for epoch in range(num_epochs):
8     #Shuffle just mixes up the dataset between epocs
9     X_train, y_train = shuffle(train_x, train_y)
10    # Mini batch learning
11    for i in range(batch_no):
12        start = i * batch_size
13        end = start + batch_size
14        inputs = Variable(torch.FloatTensor(X_train[start:end]))
15        labels = Variable(torch.FloatTensor(y_train[start:end]))
16        # zero the parameter gradients
17        optimizer.zero_grad()
18
19        # forward + backward + optimize
20        outputs = model(inputs)
21        #print("outputs",outputs)
22        #print("outputs",outputs,outputs.shape,"labels",labels, labels.shape)
23        loss = criterion(outputs, torch.unsqueeze(labels,dim=1))
24        loss.backward()
25        optimizer.step()
26
27        # print statistics
28        running_loss += loss.item()
29        losses.append(running_loss)
30
31    print('Epoch {}'.format(epoch+1), "loss: ",running_loss)
32    running_loss = 0.0
33
34 plt.plot(losses)
35 plt.savefig('loss_1.jpg')
36 plt.show()
37

```

```

Epoch 971 loss: 0.18322661785441596
Epoch 972 loss: 0.10760269820821122
Epoch 973 loss: 0.1769741404123124
Epoch 974 loss: 0.26035080274050415
Epoch 975 loss: 0.24881812533294578
Epoch 976 loss: 0.17376035023607983
Epoch 977 loss: 0.7506302666920419
Epoch 978 loss: 0.5750932950868446
Epoch 979 loss: 0.49313112707568507
Epoch 980 loss: 0.41810559765417565
Epoch 981 loss: 0.7280912714759324
Epoch 982 loss: 0.3019555547089112
Epoch 983 loss: 0.40656463149935007
Epoch 984 loss: 0.24431196725709015
Epoch 985 loss: 0.19043419955960417
Epoch 986 loss: 0.2741489303598428
Epoch 987 loss: 0.19289281975216
Epoch 988 loss: 0.09748574052900949
Epoch 989 loss: 0.06079437142579991
Epoch 990 loss: 0.18782002945954446
Epoch 991 loss: 0.22536235321331333
Epoch 992 loss: 0.10690255970348517
Epoch 993 loss: 0.09410200575894123
Epoch 994 loss: 0.10763223485992057
Epoch 995 loss: 0.15392030361999787
Epoch 996 loss: 0.26773202342064906
Epoch 997 loss: 0.15198937654986366
Epoch 998 loss: 0.07009553032912663
Epoch 999 loss: 0.08804050797061791
Epoch 1000 loss: 0.04363294974427845

```



결과 및 결론

(Conclusion)

92%의 정확도를 보임

DAICON 커뮤니티 대회 교육 랭킹 더보기

항공사 고객 만족도 예측 경진대회
데이콘 베이직 Basic | 정확도 | accuracy
₩ 상금 : 참가시 최소 50 XP, 특별상 데이콘 후드
🕒 2022.02.07 ~ 2022.02.18 17:59 + Google Calendar
👤 566명 📅 D-1

대회안내 데이터 코드 공유 토크 리더보드 팀 제출

PUBLIC RANKING CHART 순위기준

● WINNER ● 1% ● 4% ● 10%

| # | 팀 | 팀 멤버 | 점수 | 제출수 | 등록일 |
|-----|------|------|------|-----|-------|
| 156 | 김선우우 | 김선 | 0.92 | 4 | 3시간 전 |