



# 온라인 채널 제품 판매량 예측 AI 온라인 해커톤

Team KyleYeo

여훈기  
최재무  
정찬영  
신건우

# Contents

01

**Data  
Interpolation**

02

**Feature  
Engineering**

03

**Min-Max Scaling  
Log Transformation**

04

**Model**

## Setting

### Import

```
1 import random
2 import os
3 import multiprocessing
4 from tqdm.auto import tqdm
5 import warnings
6 warnings.filterwarnings('ignore')
7
8 import pandas as pd
9 import numpy as np
10 import matplotlib.pyplot as plt
11 from scipy import spatial
12
13 from sklearn.preprocessing import LabelEncoder
14 from sklearn.preprocessing import MinMaxScaler
15
16 import torch
17 import torch.nn as nn
18 import torch.optim as optim
19 import torch.nn.functional as F
20 from torch.utils.data import Dataset, DataLoader
```

### Google Drive Mount

```
1 from google.colab import drive
2
3 drive.mount('/content/drive')
```

### Hyperparameter Setting

```
1 CFG = {
2     'TRAIN_WINDOW_SIZE':60, # 60일치로 학습
3     'PREDICT_SIZE':21, # 21일치 예측
4     'EPOCHS':20,
5     'LEARNING_RATE':0.000075,
6     'BATCH_SIZE':1024,
7     'SEED':777
8 }
```

Colab python version : 3.10.12

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

### Versions

- pandas - 1.5.3
- numpy - 1.23.5
- matplotlib - 3.7.1
- tqdm - 4.66.1
- sklearn - 1.2.2
- scipy - 1.10.1
- torch - 2.0.1+cu118 ver.

### OS

- Google Colab Pro+
- GPU - V100
- 고용량 RAM 사용

# Data Interpolation

The chart displays sales data for five categories from January 2022 to April 2023. The Y-axis represents sales volume, ranging from 0 to 300,000. The X-axis represents the date. Category B002-C001-0002 (orange) is the dominant performer, with sales fluctuating between approximately 90,000 and 290,000 until late 2022, after which it drops sharply to near zero. Category B002-C001-0001 (blue) shows moderate sales, peaking at around 115,000 in early 2022 and then declining to near zero by late 2022. Categories B002-C001-0003 (green), B002-C001-0004 (red), and B002-C001-0005 (purple) maintain low sales levels throughout the period, generally below 20,000. A significant drop in sales for all categories is observed starting around late 2022, with a slight recovery in early 2023.

Date	B002-C001-0001	B002-C001-0002	B002-C001-0003	B002-C001-0004	B002-C001-0005
2022-01-01	80,000	210,000	5,000	2,000	15,000
2022-02-01	90,000	260,000	5,000	2,000	15,000
2022-03-01	115,000	290,000	5,000	2,000	15,000
2022-04-01	100,000	210,000	5,000	2,000	15,000
2022-05-01	30,000	90,000	5,000	2,000	10,000
2022-06-01	90,000	280,000	5,000	2,000	15,000
2022-07-01	65,000	170,000	5,000	2,000	10,000
2022-08-01	85,000	260,000	5,000	2,000	15,000
2022-09-01	65,000	210,000	5,000	2,000	15,000
2022-10-01	95,000	255,000	5,000	2,000	15,000
2022-11-01	70,000	165,000	5,000	2,000	10,000
2022-12-01	65,000	150,000	5,000	2,000	10,000
2023-01-01	80,000	220,000	5,000	2,000	10,000
2023-02-01	75,000	135,000	5,000	2,000	10,000
2023-03-01	10,000	10,000	5,000	2,000	10,000
2023-04-01	45,000	100,000	5,000	2,000	10,000
2023-05-01	35,000	70,000	5,000	2,000	10,000
2023-06-01	55,000	95,000	5,000	2,000	10,000

## Data Interpolation

EDA를 진행하는 도중 결측치로 의심되는 구간을 발견했다. 내부 토의로 보간하기로 결정하고, 여러 방식으로 보간을 시도한 결과, Moving Average + previous & next values를 이용하여 Simplicity하게 보간된 데이터가 가장 좋은 성능을 보였다.

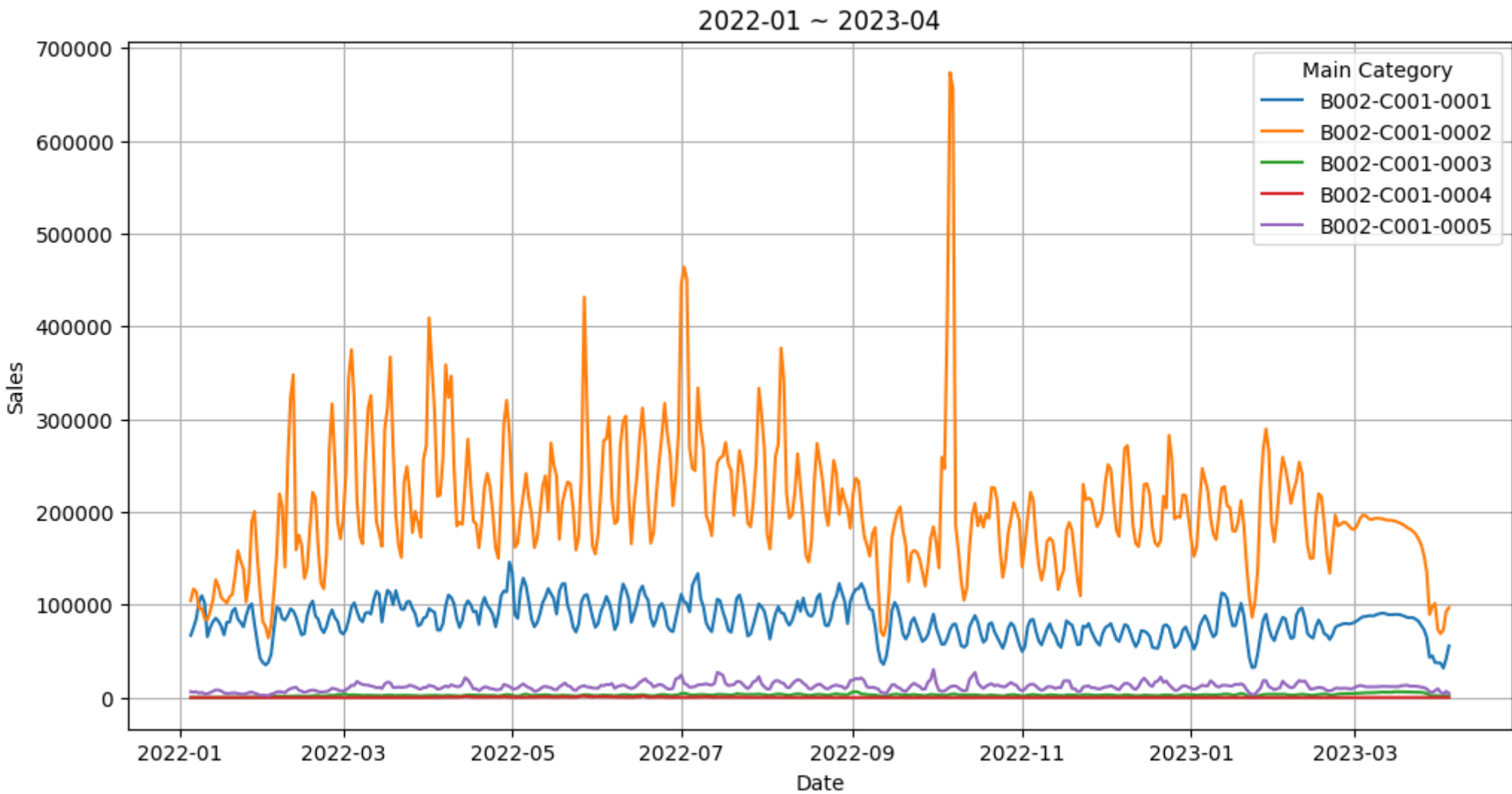
### interpolation

```
1 for i in range(0, 15890):
2     a, b = 414, 424 # 결측 구간 이전 10일 구간
3     ee = 34 # 결측 일수
4
5     z1 = train_df.iloc[i, -7:].sum() # 결측 구간 이후 7일동안의 값 더함
6     z1 = round(z1 / 7) # 7일동안의 값 평균 계산
7
8     for j in range(1, 35): # 결측 구간 보간 진행
9         z = train_df.iloc[i, a:b].sum() # 결측 구간 이전 10일동안의 값 더함.
10        z = round(z / 10) # 10일동안의 값 평균 계산
11        sum = round(z - ((z - z1) / ee)) # 결측 이전 시점의 평균값과 결측 이후 시점의 평균값을 뺀 상태에서 남은 보간 일수만큼 나누어줌. 이후 결측 이전 시점의 평균값에 빼줌
12        if sum <= 0: # 0일 경우 넘어감.
13            continue
14        else:
15            train_df.iloc[i, b] += sum # 결측 구간에 계산된 값 삽입
16            a += 1
17            b += 1
18            ee -= 1
19
20 train_df.tail(20)
```





데이터 보간 후 그래프



# Feature Engineering



## Feature Engineering

추가한 feature는 다음과 같다.

월별 판매량 기반으로 브랜드 충성도를 정의

brand\_keyword\_cnt 데이터 활용

**brand**

마지막 3주 간의  
관련 키워드 언급 횟수를  
정규화하여 평균

**Popularity1,2**

Approach1:  
월별 평균 판매량이 높을수록  
브랜드의 충성도는 높을 것이다.

Approach2:  
월별 평균 판매량이 높고 일정할수록  
브랜드의 충성도는 높을 것이다.

**Popularity3,4**

베이스라인 train 기간인  
90일 기준으로 평균을 계산

상관관계 분석을 통해,  
Popularity1,4를 최종 선정

## ‘brand’ feature

brand\_keyword\_cnt(브랜드 키워드 언급량) 데이터와 당일~몇일 후의 판매량과 상관관계가 일부 있음을 확인했다. 여러 Test를 진행하였고, 가장 최근의 Trend에 비중을 두는 것이 정확도 향상에 도움이 될 것이라 판단하여 예측 기간 직전 3주간의 평균을 feature로 사용했다.

```
1 # brand_keyword_cnt.csv 마지막 3주 값만 가지는 DataFrame으로 변환
2
3 # 원하는 날짜 범위
4 start_date = '2023-03-15'
5 end_date = '2023-04-04'
6
7 # DataFrame 변환
8 df = pd.DataFrame(brand_df)
9 date_columns = pd.date_range(start_date, end_date)
10 df_selected = df[['브랜드'] + date_columns.strftime('%Y-%m-%d').tolist()]
11
12 # DataFrame 생성
13 df = pd.DataFrame(df_selected)
14
15 # '2023-03-05'부터 마지막 컬럼까지의 데이터 선택
16 date_columns = df.columns[1:]
17
18 # 각 행별로 해당 날짜 컬럼들의 값을 정규화한 뒤 평균 계산
19 normalized_values = df[date_columns].apply(lambda row: MinMaxScaler().fit_transform(row.values.reshape(-1, 1)).mean(), axis=1)
20
21 # 정규화된 평균 값을 새로운 컬럼에 추가하되, NaN 값을 0으로 변경
22 df['brand'] = normalized_values.fillna(0)
23
24 df['brand']
```



‘brand’ feature

Train dataframe에 추출된 ‘brand’ feature를 추가한 모습

ID		제품	대분류	중분류	소분류	브랜드	brand
0	0	B002-00001-00001	B002-C001-0002	B002-C002-0007	B002-C003-0038	B002-00001	0.562504
1	1	B002-00002-00001	B002-C001-0003	B002-C002-0008	B002-C003-0044	B002-00002	0.301587
2	2	B002-00002-00002	B002-C001-0003	B002-C002-0008	B002-C003-0044	B002-00002	0.301587
3	3	B002-00002-00003	B002-C001-0003	B002-C002-0008	B002-C003-0044	B002-00002	0.301587
4	4	B002-00003-00001	B002-C001-0001	B002-C002-0001	B002-C003-0003	B002-00003	0.316015

## ‘popularity1,4’ feature

```
1 # 브랜드 충성도 feature engineering
2
3 # - Approach1: 월별 평균 판매량이 높을수록 브랜드의 충성도는 높을 것이다.
4 #   - 충성도 측정 방법: 브랜드별 월별 판매량 및 평균 판매량 계산 후 정규화 -> 값이 높을수록 충성도가 높은 브랜드임을 가정
5
6 # - Approach2: 월별 평균 판매량이 높고 일정할수록 브랜드의 충성도는 높을 것이다.
7 #   - 충성도 측정 방법: 평균 월별 판매량과 CV (Coefficient of Variation) 기반으로 각각 가중치 적용 후 weighted mean 계산 -> 값이 높을수록 충성도가 높은 브랜드임을 가정
8
9 class brand_popularity:
10     def __init__(self, df: pd.DataFrame, w1, w2, target_df: pd.DataFrame):
11         """
12         df: 브랜드별 월별 판매량이 산출된 df,
13         w1: 월별 판매량에 적용할 가중치
14         w2: 월별 판매량의 일정함(표준편차)에 적용할 가중치
15         target_df: 계산된 feature를 mapping할 df
16         """
17         self.df = df
18         self.w1 = w1
19         self.w2 = w2
20         self.target_df = target_df
21
22     def sales_based_popularity(self):
23         """
24         브랜드별 월별 평균 판매량을 활용한 충성도 점수
25         """
26         average_monthly_sales = self.df.mean()
27         normalized_popularity = average_monthly_sales / average_monthly_sales.max()
28         return normalized_popularity
```



## ‘popularity1,4’ feature

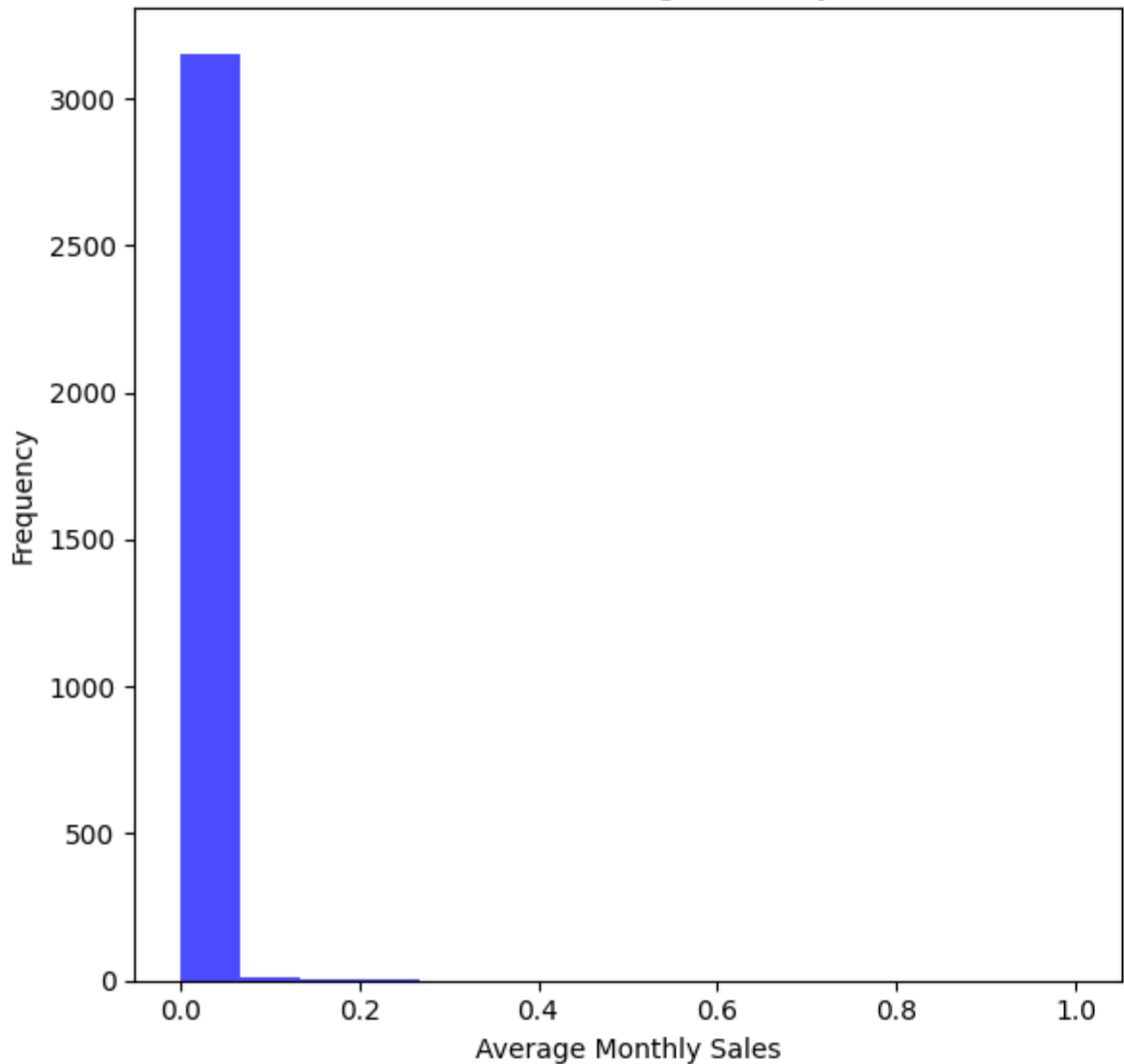
```
30 def mixed_based_popularity(self):
31     ...
32     월별 판매량과 판매량의 표준편차 (cv)를 활용한 충성도 점수
33     ...
34     average_monthly_sales = self.df.mean()
35     normalized_popularity = average_monthly_sales / average_monthly_sales.max()
36     cv = self.df.std() / self.df.mean()
37     normalized_cv = cv / cv.max()
38
39     combined_score = (self.w1 * normalized_popularity) + (self.w2 * (1 - normalized_cv))
40     return combined_score
```



‘popularity1,4’ feature

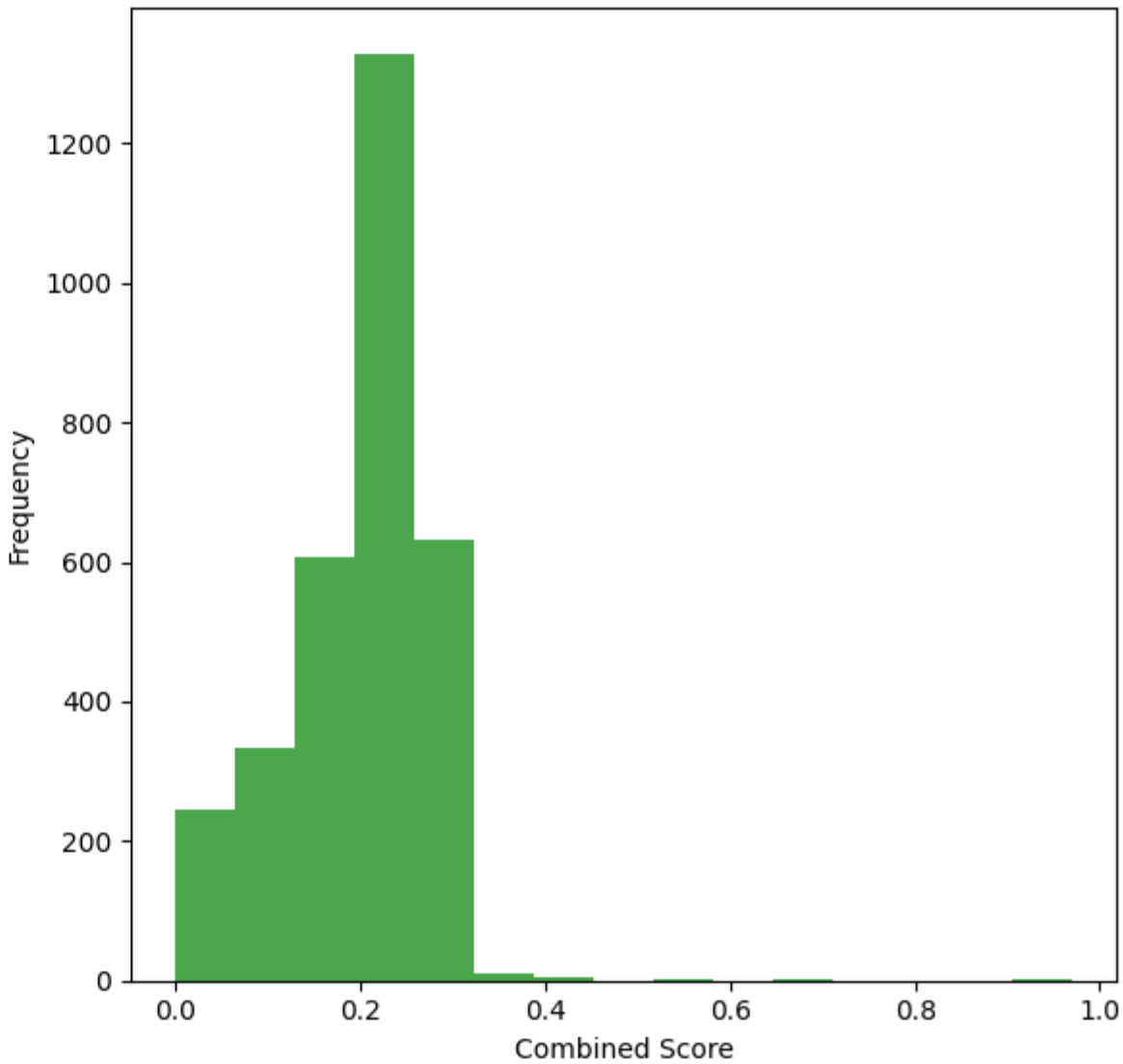
popularity1

Distribution of Average Monthly Sales



popularity2

Distribution of Combined Score

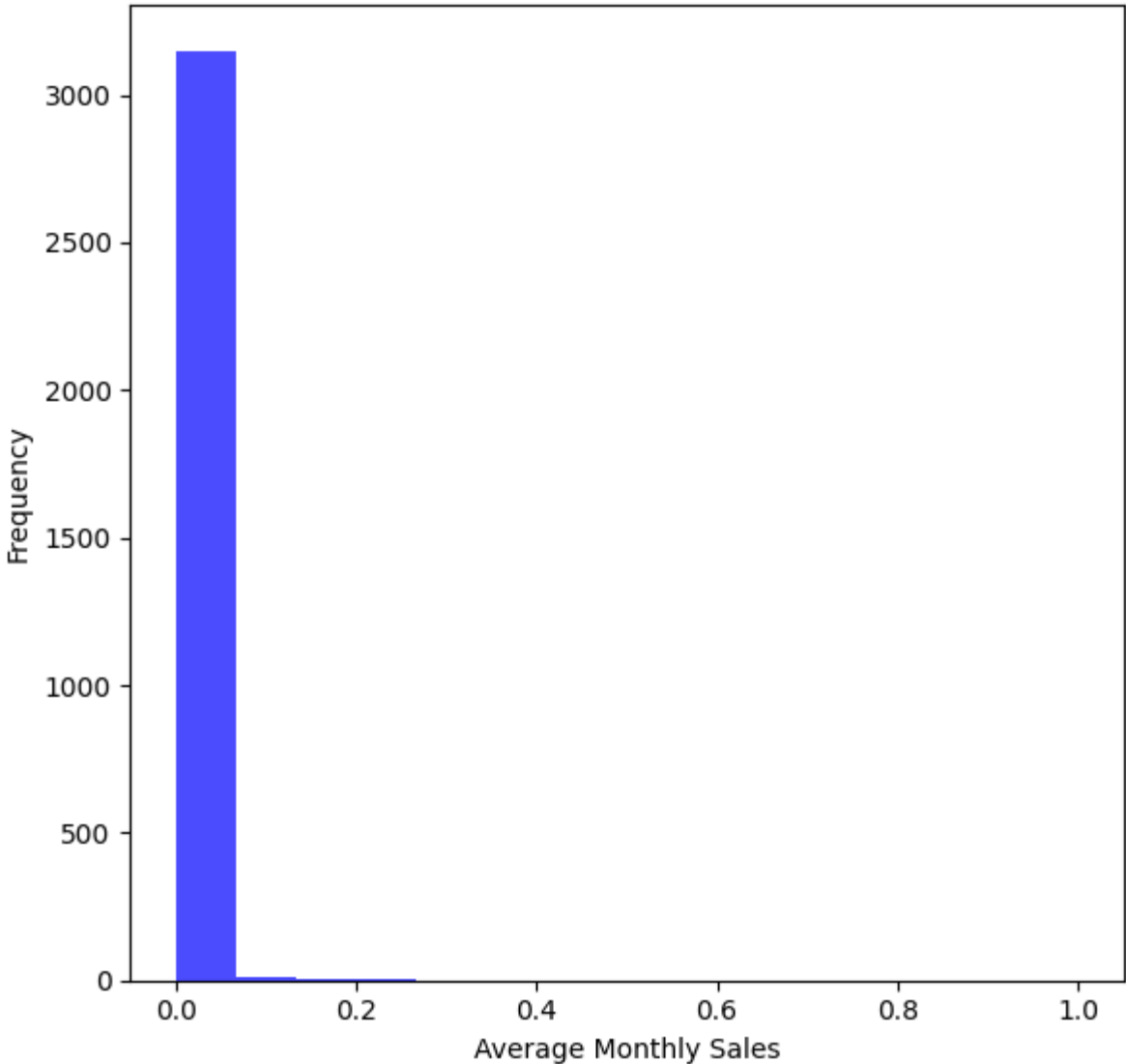




‘popularity1,4’ feature

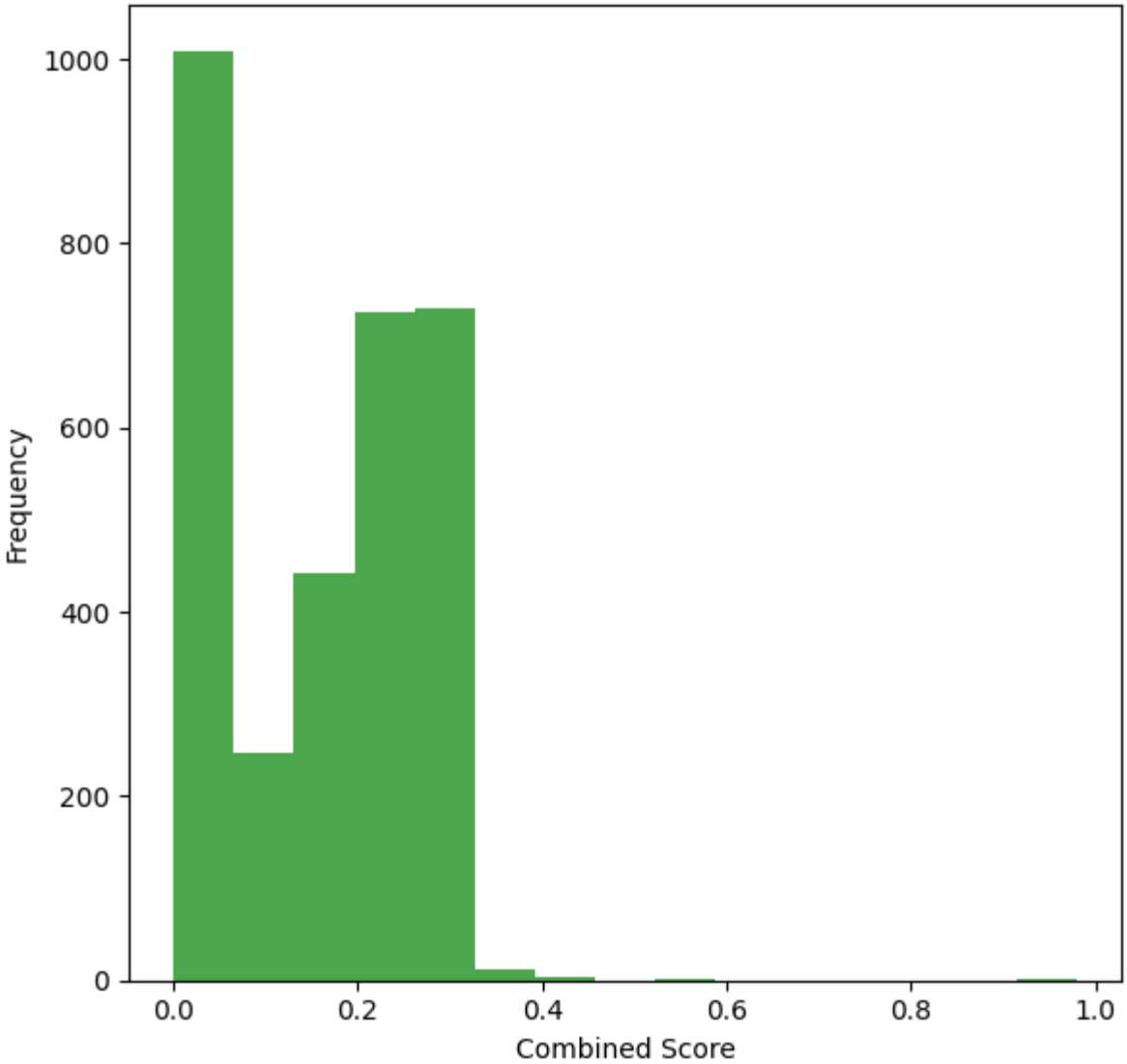
popularity3

Distribution of Average Monthly Sales



popularity4

Distribution of Combined Score





## ‘popularity1,4’ feature

```
1 # Pearson correlation matrix
2 df_features.corr()
```

	popularity1	popularity2	popularity3	popularity4	total_sales	avg_sales
popularity1	1.000000	0.854671	0.975166	0.692979	0.243066	0.243066
popularity2	0.854671	1.000000	0.839259	0.858462	0.214701	0.214701
popularity3	0.975166	0.839259	1.000000	0.707194	0.204916	0.204916
popularity4	0.692979	0.858462	0.707194	1.000000	0.159251	0.159251
total_sales	0.243066	0.214701	0.204916	0.159251	1.000000	1.000000
avg_sales	0.243066	0.214701	0.204916	0.159251	1.000000	1.000000

```
1 # 비선형 관계를 파악하기 위해 각 feature와 브랜드별 총 판매량 사이의 Distance correlation coefficients 계산합니다
2 display(spatial.distance.correlation(df_features['popularity1'], df_features['total_sales']))
3 display(spatial.distance.correlation(df_features['popularity2'], df_features['total_sales']))
4 display(spatial.distance.correlation(df_features['popularity3'], df_features['total_sales']))
5 display(spatial.distance.correlation(df_features['popularity4'], df_features['total_sales']))
```

```
0.7569341821112185
0.7852989927187264
0.7950839016496316
0.8407490358050872
```

```
1 # 비선형 관계를 파악하기 위해 각 feature와 브랜드별 평균 판매량 사이의 Distance correlation coefficients 계산합니다
2 display(spatial.distance.correlation(df_features['popularity1'], df_features['avg_sales']))
3 display(spatial.distance.correlation(df_features['popularity2'], df_features['avg_sales']))
4 display(spatial.distance.correlation(df_features['popularity3'], df_features['avg_sales']))
5 display(spatial.distance.correlation(df_features['popularity4'], df_features['avg_sales']))
```

```
0.7569341821112185
0.7852989927187264
0.7950839016496316
0.8407490358050872
```

### Popularity1

● 판매량과 그나마 가장 강한 선형적인 관계를 나타냄  
(0.24로 의미부여하기 어려운 상관계수이긴 하지만)  
또 강한 비선형 관계를 나타냄 (0.76)

### Popularity4

● 판매량과 가장 강한 비선형 관계를 나타냄 (0.84)



**Popularity1,4를 최종 선정**

**Min-Max Scaling**  
**Log Transformation**

## Min-Max Scaling

베이스라인의 Min-max scaling을 그대로 적용했다.

```
1 # 숫자형 변수들의 min-max scaling을 수행하는 코드입니다.
2 numeric_cols = train_df.columns[7:]
3 # 각 column의 min 및 max 계산
4 min_values = train_df[numeric_cols].min(axis=1)
5 max_values = train_df[numeric_cols].max(axis=1)
6 # 각 행의 범위(max-min)를 계산하고, 범위가 0인 경우 1로 대체
7 ranges = max_values - min_values
8 ranges[ranges == 0] = 1
9 # min-max scaling 수행
10 train_df[numeric_cols] = (train_df[numeric_cols].subtract(min_values, axis=0)).div(ranges, axis=0)
11 # max와 min 값을 dictionary 형태로 저장
12 scale_min_dict = min_values.to_dict()
13 scale_max_dict = max_values.to_dict()
```

## Log transformation

학습 데이터의 분산 안정화를 위해 로그 변환을 적용했다.  
변환된 학습 데이터로 인해 모델의 예측 성능이 좋아졌다.

```
1 # Variance 안정화를 위해 로그 변환을 적용합니다.
2
3 def log_transformation(df: pd.DataFrame, start_date, offset = float):
4     '''
5     정상성 변환을 위해 로그 변환을 적용합니다.
6     0이 많은 데이터 특성 때문에 일반적인 로그 변환보다는, 전체 데이터에 offset을 더하고 로그 변환을 진행합니다.
7     '''
8     df_log = df.copy()
9     df_log = df_log.iloc[:, df.columns.get_loc(start_date):] + offset
10    df_log = np.log(df_log)
11    df.iloc[:, df.columns.get_loc('2022-01-01'):] = df_log
12
13    return df
14
15 train_data = log_transformation(train_df, '2022-01-01', 1)
```

## 학습/검정 데이터 생성

베이스라인의 학습데이터 생성 함수에서 step size를 2로 추가해 RAM을 절약했다.  
결과 자체도 step size를 부여한 것이 성능이 더 좋았다.

```
1 def make_train_data(data, train_size=CFG['TRAIN_WINDOW_SIZE'], predict_size=CFG['PREDICT_SIZE']):
2     STEP_SIZE = 2
3
4     num_rows = len(data)
5     window_size = train_size + predict_size
6     adjusted_size = (len(data.columns) - window_size + 1) // STEP_SIZE
7
8     input_data = np.empty((num_rows * adjusted_size, train_size, len(data.iloc[0, :7]) + 1))
9     target_data = np.empty((num_rows * adjusted_size, predict_size))
10
11     for i in tqdm(range(num_rows)):
12         encode_info = np.array(data.iloc[i, :7])
13         sales_data = np.array(data.iloc[i, 7:])
14
15         for j in range(0, len(sales_data) - window_size + 1, STEP_SIZE):
16             window = sales_data[j: j + window_size]
17             temp_data = np.column_stack((np.tile(encode_info, (train_size, 1)), window[:train_size]))
18             input_data[i * adjusted_size + j // STEP_SIZE] = temp_data
19             target_data[i * adjusted_size + j // STEP_SIZE] = window[train_size:]
20
21     return input_data, target_data
```

## 학습/검정 데이터 생성

학습/검정 데이터 분할과 DataLoader 또한 베이스라인의 코드를 그대로 사용했다.  
CV 또한 적용해보았지만, CV를 적용하지 않은 것이 Public score가 더 높았다.

```
1 # Train / Validation Split
2 data_len = len(train_input)
3 val_input = train_input[-int(data_len*0.2):]
4 val_target = train_target[-int(data_len*0.2):]
5 train_input = train_input[:-int(data_len*0.2)]
6 train_target = train_target[:-int(data_len*0.2)]
```

만들어진 데이터의 shape은 다음과 같다.

```
1 train_input.shape, train_target.shape, val_input.shape, val_target.shape, test_input.shape
```

```
((2453416, 60, 8),
 (2453416, 21),
 (613354, 60, 8),
 (613354, 21),
 (15890, 60, 8))
```

**Model**



## 모델 선언

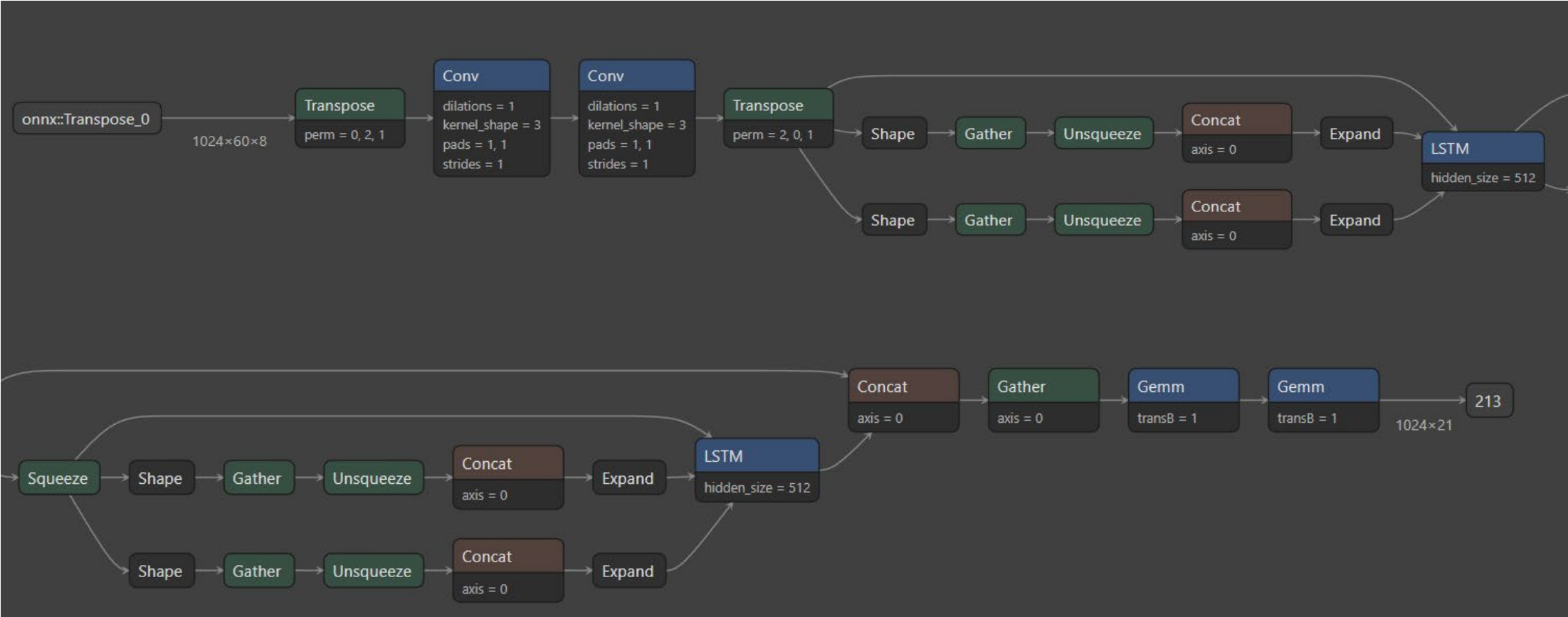
다음과 같이 Conv1d layer와 LSTM layer를 결합한 모델을 만들었다.

```
4 class Conv1dLSTMModel(nn.Module):
5     def __init__(self, input_size=8, hidden_size=512, output_size=CFG['PREDICT_SIZE']):
6         super(Conv1dLSTMModel, self).__init__()
7         self.conv1d_1 = nn.Conv1d(in_channels=input_size,
8                                     out_channels=16,
9                                     kernel_size=3,
10                                    stride=1,
11                                    padding=1)
12         self.conv1d_2 = nn.Conv1d(in_channels=16,
13                                     out_channels=32,
14                                     kernel_size=3,
15                                     stride=1,
16                                     padding=1)
17         self.lstm = nn.LSTM(input_size=32,
18                             hidden_size=hidden_size,
19                             num_layers=2,
20                             bias=True,
21                             bidirectional=False,
22                             batch_first=True)
23         self.dropout = nn.Dropout(0.1)
24
25         self.fc_layer1 = nn.Linear(512, 32)
26         self.fc_layer2 = nn.Linear(32, output_size)
```

```
28     def forward(self, x):
29         x = x.transpose(1, 2)
30
31         x = self.conv1d_1(x)
32
33         x = self.conv1d_2(x)
34
35         x = x.transpose(1, 2)
36
37         self.lstm.flatten_parameters()
38
39         _, (hidden, _) = self.lstm(x)
40
41         x = hidden[-1]
42
43         x = self.dropout(x)
44
45         x = self.fc_layer1(x)
46         x = self.fc_layer2(x)
47
48         return x
```

모델 선언

Netron으로 시각화한 모델 구조다.



## 모델 학습

```
1 def train_with_early_stopping(model, optimizer, train_loader, val_loader, device, patience=3):
2     model.to(device)
3     criterion = nn.MSELoss().to(device)
4     best_loss = float('inf')
5     best_model = None
6     no_improvement_count = 0 # 변수 초기화
7
8     for epoch in range(1, CFG['EPOCHS']+1):
9         model.train()
10        train_loss = []
11
12        for X, Y in tqdm(iter(train_loader)):
13            X = X.to(device)
14            Y = Y.to(device)
15
16            optimizer.zero_grad()
17
18            output = model(X)
19            loss = criterion(output, Y)
20
21            loss.backward()
22            optimizer.step()
23
24            train_loss.append(loss.item())
25
26        val_loss = validation(model, val_loader, criterion, device)
27        print(f'Epoch : [{epoch}] Train Loss : [{np.mean(train_loss):.5f}] Val Loss : [{val_loss:.5f}]')
28
29        if val_loss < best_loss: # 검증 손실이 감소한 경우
30            best_loss = val_loss
31            best_model = model
32            no_improvement_count = 0 # 카운트 초기화
33            print('Model Saved')
34        else:
35            no_improvement_count += 1
36
37        if no_improvement_count >= patience: # 일정 epoch 동안 검증 손실이 향상되지 않으면 종료
38            print(f'Early stopping! No improvement for {patience} epochs.')
39            print(f'TRAIN_WINDOW_SIZE : {CFG["TRAIN_WINDOW_SIZE"]}, PREDICT_SIZE : {CFG["PREDICT_SIZE"]},
40            break
41
42    return best_model
```

```
1 def validation(model, val_loader, criterion, device):
2     model.eval()
3     val_loss = []
4
5     with torch.no_grad():
6         for X, Y in tqdm(iter(val_loader)):
7             X = X.to(device)
8             Y = Y.to(device)
9
10            output = model(X)
11            loss = criterion(output, Y)
12
13            val_loss.append(loss.item())
14    return np.mean(val_loss)
```

학습 시간 절약을 위해 early stop을 적용하였고,  
Patience는 3으로 설정했다.

## 모델 예측, 역변환

베이스라인 코드를 사용하여 추론을 진행했다.  
예측값을 역변환하고 반올림하여 후처리했다.

```
1 # inverse log transformation
2 offset = 1
3 pred = np.exp(pred) - offset
4
5 # 추론 결과를 inverse scaling
6 for idx in tqdm(range(len(pred))):
7     pred.iloc[idx, :] = pred.iloc[idx, :] * (scale_max_dict[idx] - scale_min_dict[idx]) + scale_min_dict[idx]
8
9 # 결과 후처리
10 pred = np.round(pred, 0).astype(int)
```