


```

inflating: 외부데이터/USD_MYR 내역.csv
외부데이터/USD_NOK 내역.csv: mismatching "local" filename ( ЪЩҀыѰАыН ЪЭҀэДЉ/USD_N
continuing with "central" filename version
inflating: 외부데이터/USD_NOK 내역.csv
외부데이터/USD_PEN 내역.csv: mismatching "local" filename ( ЪЩҀыѰАыН ЪЭҀэДЉ/USD_P
continuing with "central" filename version
inflating: 외부데이터/USD_PEN 내역.csv
외부데이터/USD_THB 내역.csv: mismatching "local" filename ( ЪЩҀыѰАыН ЪЭҀэДЉ/USD_T
continuing with "central" filename version
inflating: 외부데이터/USD_THB 내역.csv
외부데이터/USD_VND 내역.csv: mismatching "local" filename ( ЪЩҀыѰАыН ЪЭҀэДЉ/USD_V
continuing with "central" filename version
inflating: 외부데이터/USD_VND 내역.csv
외부데이터/국제유가2021-09-09.csv: mismatching "local" filename ( ЪЩҀыѰАыН ЪЭҀэДЉ
continuing with "central" filename version
inflating: 외부데이터/국제유가2021-09-09.csv
제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_문제데이터.xlsx: mismatching "
continuing with "central" filename version
inflating: 제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_문제데이터.xlsx
제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_자율평가데이터.xlsx: mismatchi
continuing with "central" filename version
inflating: 제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_자율평가데이터.xls:
제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_평가데이터_update_210831.xlsx:
continuing with "central" filename version
inflating: 제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_평가데이터_update_

```

1 ## 한글폰트 깨짐현상 해결

2 !sudo apt-get install -y fonts-nanum

3 !sudo fc-cache -fv

4 !rm ~/.cache/matplotlib -rf

5

6 !pip install sktime ##시계열 train, test 분리할때 사용

```

debconf: unable to initialize frontend: Dialog
debconf: (No usable dialog-like program is installed, so the dialog based frontend cannot
debconf: falling back to frontend: Readline
debconf: unable to initialize frontend: Readline
debconf: (This frontend requires a controlling tty.)
debconf: falling back to frontend: Teletype
dpkg-preconfigure: unable to re-open stdin:
Selecting previously unselected package fonts-nanum.
(Reading database ... 148492 files and directories currently installed.)
Preparing to unpack .../fonts-nanum_20170925-1_all.deb ...
Unpacking fonts-nanum (20170925-1) ...
Setting up fonts-nanum (20170925-1) ...
Processing triggers for fontconfig (2.12.6-0ubuntu2) ...
/usr/share/fonts: caching, new cache contents: 0 fonts, 1 dirs
/usr/share/fonts/truetype: caching, new cache contents: 0 fonts, 3 dirs

/usr/share/fonts/truetype/humor-sans: caching, new cache contents: 1 fonts, 0 dirs
/usr/share/fonts/truetype/liberation: caching, new cache contents: 16 fonts, 0 dirs
/usr/share/fonts/truetype/nanum: caching, new cache contents: 10 fonts, 0 dirs
/usr/local/share/fonts: caching, new cache contents: 0 fonts, 0 dirs
/root/.local/share/fonts: skipping, no such directory
/root/.fonts: skipping, no such directory
/var/cache/fontconfig: cleaning cache directory
/root/.cache/fontconfig: not cleaning non-existent cache directory
/root/.fontconfig: not cleaning non-existent cache directory
fc-cache: succeeded
Collecting sktime

```

```

Downloading sktime-0.7.0-cp37-cp37m-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (5.8
|████████████████████████████████████████| 5.8 MB 5.1 MB/s
Requirement already satisfied: pandas>=1.1.0 in /usr/local/lib/python3.7/dist-packages (fr
Requirement already satisfied: wheel in /usr/local/lib/python3.7/dist-packages (from sktim
Requirement already satisfied: numpy>=1.19.0 in /usr/local/lib/python3.7/dist-packages (fr
Collecting scikit-learn>=0.24.0
  Downloading scikit_learn-0.24.2-cp37-cp37m-manylinux2010_x86_64.whl (22.3 MB)
    |████████████████████████████████████████| 22.3 MB 1.5 MB/s
Requirement already satisfied: numba>=0.50 in /usr/local/lib/python3.7/dist-packages (from
Collecting statsmodels>=0.12.1
  Downloading statsmodels-0.12.2-cp37-cp37m-manylinux1_x86_64.whl (9.5 MB)
    |████████████████████████████████████████| 9.5 MB 54.5 MB/s
Requirement already satisfied: setuptools in /usr/local/lib/python3.7/dist-packages (from
Requirement already satisfied: llvmlite<0.35,>=0.34.0.dev0 in /usr/local/lib/python3.7/dis
Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-pac
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.7/dist-packages (fr
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.7/dist-packages (from py
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.7/dist-packages (fr
Requirement already satisfied: scipy>=0.19.1 in /usr/local/lib/python3.7/dist-packages (fr
Collecting threadpoolctl>=2.0.0
  Downloading threadpoolctl-2.2.0-py3-none-any.whl (12 kB)
Requirement already satisfied: patsy>=0.5 in /usr/local/lib/python3.7/dist-packages (from
Installing collected packages: threadpoolctl, statsmodels, scikit-learn, sktime
  Attempting uninstall: statsmodels
    Found existing installation: statsmodels 0.10.2
    Uninstalling statsmodels-0.10.2:
      Successfully uninstalled statsmodels-0.10.2
  Attempting uninstall: scikit-learn
    Found existing installation: scikit-learn 0.22.2.post1
    Uninstalling scikit-learn-0.22.2.post1:
      Successfully uninstalled scikit-learn-0.22.2.post1
Successfully installed scikit-learn-0.24.2 sktime-0.7.0 statsmodels-0.12.2 threadpoolctl-2

```

런타임 재시작하기

▼ 2. 모듈 및 데이터불러오기

```

1 ## 기본적인 모듈
2 import pandas as pd
3 import warnings ; warnings.filterwarnings(action='ignore')
4 import numpy as np
5 from datetime import datetime
6 from tqdm.notebook import tqdm ## 진행상황 표시
7
8 ## 시각화 및 폰트 지정할때 사용
9 import matplotlib.pyplot as plt ; plt.rc('font', family='NanumBarunGothic',size=20) ## 기본 시크
10 import seaborn as sns ## 기본 시각화
11 from sklearn.metrics import silhouette_samples, silhouette_score # 실루엣 계수 시각화
12 from sklearn.preprocessing import LabelEncoder ## 클러스터링 할때 사용
13 from sklearn.cluster import KMeans ## 클러스터링
14
15 ## 모델링 할때 사용

```

```

16 from sktime.forecasting.model_selection import temporal_train_test_split ## 시계열 train, test
17 import xgboost as xgb ## 클러스터링 예측할때 사용
18 from xgboost import XGBRegressor ## 가격 예측할때 사용
19 from sklearn.metrics import mean_squared_error ## mse 계산할때 사용

1 ## 데이터 경로 지정
2 TRAIN_PATH_2015_2019 = './제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_문제데이
3 TRAIN_PATH_2020 = './제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_자율평가데이
4 TARGET_PATH_2021 = './제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_평가데이터_u
5
6 ## 2015년 ~ 2020년 학습 데이터 불러오기
7 train_full_data=pd.DataFrame()
8 train_data_2015_2019=pd.read_excel(TRAIN_PATH_2015_2019)
9 train_data_2020=pd.read_excel(TRAIN_PATH_2020)
10
11 ## 데이터 연결
12 train_full_data=pd.concat([train_data_2015_2019,train_data_2020])
13 train_full_data.reset_index(drop=True,inplace=True)

```

▼ 3. 제공 데이터 전처리

실루엣 계수 및 시각화를 이용해서 상세 어종별 최적 군집 선택

```

1 ## 상세어종이 오징어, 연어, 흰다리 새우인 데이터만 보기
2 train_full_data=train_full_data[train_full_data['P_NAME'].isin(['오징어','연어','흰다리새우'])]
3
4
5 ## 필요없는 변수들 제거
6 ## P_TYPE : 모두 같아서 제거함
7 ## CATEGORY_1, CATEGORY_2 : 원하는 제품의 경우 모든 내용이 같으므로 제거함
8 ## CTRY_2 : 제조국과 판매국이 같은경우가 약 97%정도이므로 국가 정보변수 제거 해도 큰 영향 없을거
9 train_full_data.drop(['P_TYPE','CTRY_2','CATEGORY_1','CATEGORY_2'],axis=1,inplace=True)
10
11
12 ## 제품이 가격에 큰 영향을 줄것이라 판단하여 제품을 따로 분류하고, 모델링을 진행할 것이기 때문에
13 train_data_squid=train_full_data[train_full_data['P_NAME']=='오징어'] ## 오징어 제품
14 train_data_salmon=train_full_data[train_full_data['P_NAME']=='연어'] ## 연어 제품
15 train_data_white_shrimp=train_full_data[train_full_data['P_NAME']=='흰다리새우'] ## 흰다리새우
16
17
18 ## 제품별 새로운 파생변수 생성 (제조국_수입용도_수입형태)
19 ## 제조국, 수입용도, 수입형태 모든 조합별 파생변수를 만들기에는 데이터에 비해 변수가 많이 생길것
20 for input_p_name in ['train_data_squid','train_data_salmon','train_data_white_shrimp']:
21     globals()[input_p_name]['ctry_purpose_import_type']=globals()[input_p_name].apply(lambda :

1 ## 실루엣 계수 사용
2 ## 실루엣 계수 - 실루엣 계수란 각 군집 간의 거리가 얼마나 효율적으로 분리돼 있는지를 나타내는 거
3 ## 시각화 해석 방법 1 - 계수가 1에 가까울수록 다른 군집과 거리가 멀다는 것을 의미
4 ## 시각화 해석 방법 2 - 평균값에 비해 계수가 떨어지는 군집이 존재하면 안됨
5
6 def visualize_silhouette(cluster_lists, X_features):

```

```

7
8     from sklearn.datasets import make_blobs
9     from sklearn.cluster import KMeans
10    from sklearn.metrics import silhouette_samples, silhouette_score
11    import matplotlib.pyplot as plt
12    import matplotlib.cm as cm
13    import math
14
15    # 입력값으로 클러스터링 갯수들을 리스트로 받아서, 각 갯수별로 클러스터링을 적용하고 실루엣 :
16    n_cols = len(cluster_lists)
17
18    # plt.subplots()으로 리스트에 기재된 클러스터링 수만큼의 sub figures를 가지는 axs 생성 Wn",
19    fig, axs = plt.subplots(figsize=(10*n_cols, 10), nrows=1, ncols=n_cols)
20
21
22    # 리스트에 기재된 클러스터링 갯수들을 차례로 iteration 수행하면서 실루엣 개수 시각화Wn",
23    for ind, n_cluster in enumerate(cluster_lists):
24
25        # KMeans 클러스터링 수행하고, 실루엣 스코어와 개별 데이터의 실루엣 값 계산. Wn",
26        clusterer = KMeans(n_clusters = n_cluster, max_iter=500, random_state=0)
27        cluster_labels = clusterer.fit_predict(X_features)
28
29        sil_avg = silhouette_score(X_features, cluster_labels)
30        sil_values = silhouette_samples(X_features, cluster_labels)
31
32        y_lower = 10
33        axs[ind].set_title('Number of Cluster : ' + str(n_cluster)+
34                           'Silhouette Score : ' + str(round(sil_avg,3)) )
35        axs[ind].set_xlabel("The silhouette coefficient values")
36        axs[ind].set_ylabel("Cluster label")
37        axs[ind].set_xlim([-0.1, 1])
38        axs[ind].set_ylim([0, len(X_features) + (n_cluster + 1) * 10])
39        axs[ind].set_yticks([]) # Clear the yaxis labels / ticksWn",
40        axs[ind].set_xticks([0, 0.2, 0.4, 0.6, 0.8, 1])
41
42    # 클러스터링 갯수별로 fill_betweenx( )형태의 막대 그래프 표현.
43        for i in range(n_cluster):
44            ith_cluster_sil_values = sil_values[cluster_labels==i]
45            ith_cluster_sil_values.sort()
46
47            size_cluster_i = ith_cluster_sil_values.shape[0]
48            y_upper = y_lower + size_cluster_i
49
50            color = cm.nipy_spectral(float(i) / n_cluster)
51            axs[ind].fill_betweenx(np.arange(y_lower, y_upper), 0, ith_cluster_sil_values,
52                                  facecolor=color, edgecolor=color, alpha=0.7)
53            axs[ind].text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))
54            y_lower = y_upper + 10
55
56        axs[ind].axvline(x=sil_avg, color="red", linestyle="--")
57
58    1 ## 실루엣 계수 사용을 위해 라벨인코딩
59    2 le = LabelEncoder()
60    3 train_data_squid['sil_squid'] = le.fit_transform(train_data_squid['ctry_purpose_import_type']).va
61    4

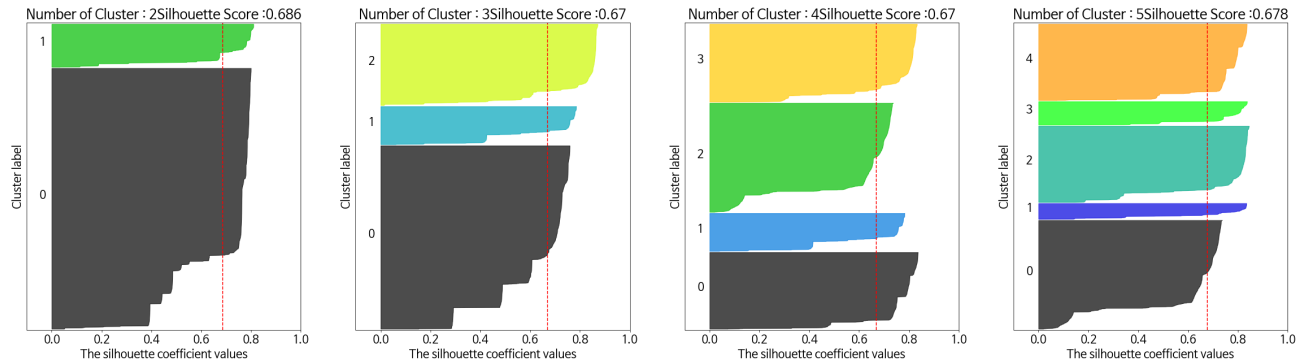
```

```

5 ## 불필요한 열 제거
6 train_data_squid_sil = train_data_squid.drop(columns=['REG_DATE', 'CTRY_1', 'P_PURPOSE', 'P_NAME'])

1 ## 실루엣 계수 시각화
2 ## 해석 결과 : 오징어의 경우 클러스터 4개가 최적이라 판단
3 visualize_silhouette([2,3,4,5],train_data_squid_sil)

```



```

1 ## 군집화
2 ## 날짜 및 파생변수(제조국_수입용도_수입형태)별 평균 금액 이용
3 def kmeans_data_squid(input_data):
4     df_C_P_I = input_data.groupby(['ctry_purpose_import_type', 'REG_DATE'])['P_PRICE'].mean().reset_index()
5     for i in range(len(df_C_P_I)):
6         # 날짜별 평균 금액에 대해 scaling 하여 군집화
7         df_C_P_I.iloc[i,1:] = (df_C_P_I.iloc[i,1:] - df_C_P_I.iloc[i,1:].mean()) / df_C_P_I.iloc[i,1:].std()
8     df_C_P_I.fillna(0, inplace=True)
9     kmeans = KMeans(n_clusters=4, random_state = 7857)
10    km_cluster = kmeans.fit_predict(df_C_P_I.iloc[:,1:])
11
12    df_C_P_I_clust = df_C_P_I.copy()
13    df_C_P_I_clust['km_cluster_C_P_I'] = km_cluster
14
15    input_data=input_data.merge(df_C_P_I_clust[['ctry_purpose_import_type', 'km_cluster_C_P_I']], on=['ctry_purpose_import_type', 'REG_DATE'], how='left')
16
17    return input_data
18

```

```

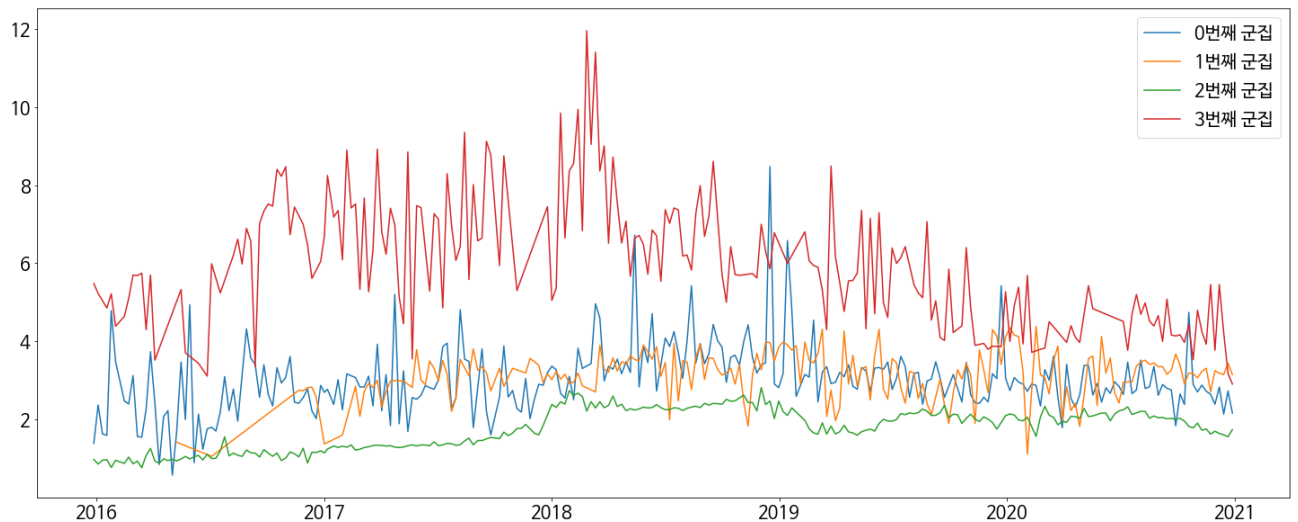
1 ## 오징어 kmeans 실행
2 train_data_squid=kmeans_data_squid(train_data_squid)
3
4 ## 오징어 시각화 준비
5 train_data_squid_cluster=train_data_squid.groupby(['REG_DATE', 'km_cluster_C_P_I'])['P_PRICE'].mean().reset_index()
6 train_data_squid_cluster.index=train_data_squid_cluster['REG_DATE']

```

```

1 ## 군집별, 날짜별 평균 금액
2 fig = plt.figure(figsize = (25, 10))
3
4 plt.plot(train_data_squid_cluster[train_data_squid_cluster['km_cluster_C_P_I']==0]['P_PRICE'], label='0번째 군집')
5 plt.plot(train_data_squid_cluster[train_data_squid_cluster['km_cluster_C_P_I']==1]['P_PRICE'], label='1번째 군집')
6 plt.plot(train_data_squid_cluster[train_data_squid_cluster['km_cluster_C_P_I']==2]['P_PRICE'], label='2번째 군집')
7 plt.plot(train_data_squid_cluster[train_data_squid_cluster['km_cluster_C_P_I']==3]['P_PRICE'], label='3번째 군집')
8
9 plt.legend(loc='best')
10 plt.show()

```



연어 실루엣 검증

```

1 ## 실루엣 계수 사용을 위해 라벨인코딩
2 train_data_salmon['sil_squid'] = le.fit_transform(train_data_salmon['ctry_purpose_import_type'])
3
4 ## 불필요한 열 제거
5 train_data_salmon_sil = train_data_salmon.drop(columns=['REG_DATE', 'CTRY_1', 'P_PURPOSE', 'P_N

```

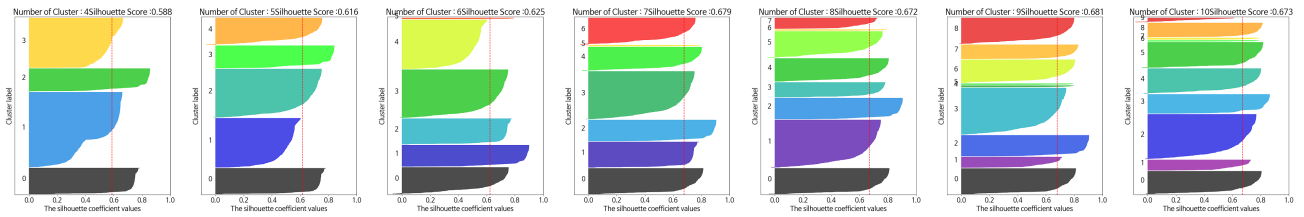
```

1 # 실루엣 계수 시각화
2 # 연어 최적 클러스터 개수는 4개
3 # 5,6 클러스터로 군집화 하였을때 실루엣 계수는 높지만 평균값(빨간선)에 못미치는 군집이 존재

```

4

5 visualize_silhouette([4,5,6,7,8,9,10],train_data_salmon_sil)



1 ## 군집화

```

2 def kmeans_data_salmon(input_data):
3     df_C_P_I = input_data.groupby(['ctry_purpose_import_type', 'REG_DATE'])['P_PRICE'].mean().reset_index()
4     for i in range(len(df_C_P_I)):
5         df_C_P_I.iloc[i,1:] = (df_C_P_I.iloc[i,1:] - df_C_P_I.iloc[i,1:].mean())/df_C_P_I.iloc[i,1:].std()
6     df_C_P_I.fillna(0,inplace=True)
7     kmeans = KMeans(n_clusters=4, random_state = 7857)
8     km_cluster = kmeans.fit_predict(df_C_P_I.iloc[:,1:])
9
10    df_C_P_I_clust = df_C_P_I.copy()
11    df_C_P_I_clust['km_cluster_C_P_I'] = km_cluster
12
13    input_data=input_data.merge(df_C_P_I_clust[['ctry_purpose_import_type', 'km_cluster_C_P_I']])
14
15    return input_data
16

```

1 ## 연어 kmeans 실행

```

2 train_data_salmon=kmeans_data_salmon(train_data_salmon)
3
4 ## 연어 시각화 준비
5 train_data_salmon_cluster=train_data_salmon.groupby(['REG_DATE', 'km_cluster_C_P_I'])['P_PRICE'].mean().reset_index()
6 train_data_salmon_cluster.index=train_data_salmon_cluster['REG_DATE']

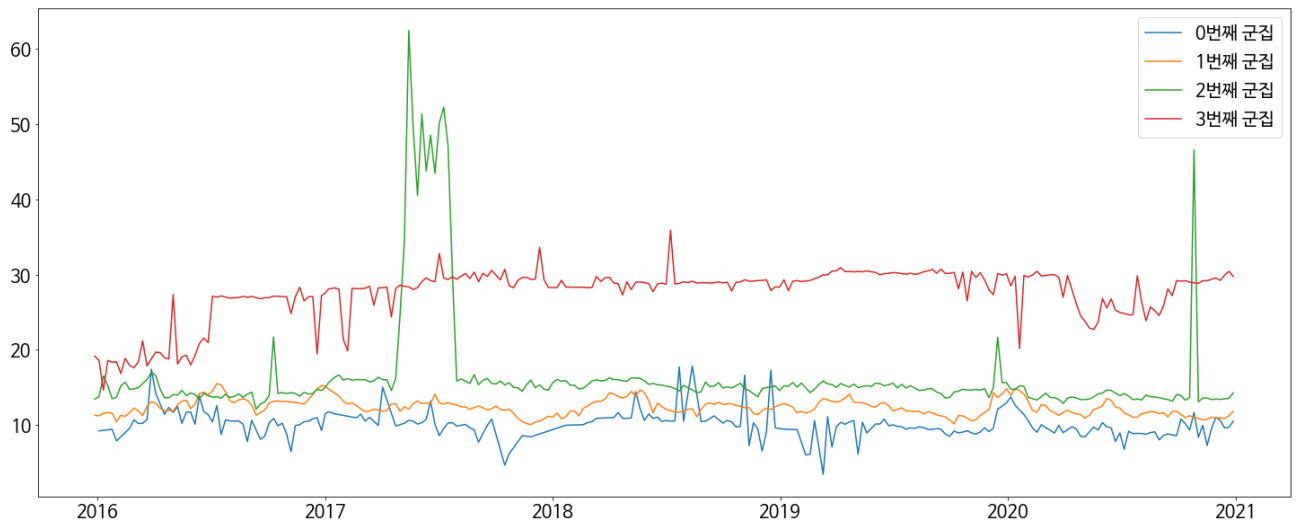
```

1 ## 시각화

```

2 fig = plt.figure(figsize = (25, 10))
3
4 plt.plot(train_data_salmon_cluster[train_data_salmon_cluster['km_cluster_C_P_I']==0]['P_PRICE'])
5 plt.plot(train_data_salmon_cluster[train_data_salmon_cluster['km_cluster_C_P_I']==1]['P_PRICE'])
6 plt.plot(train_data_salmon_cluster[train_data_salmon_cluster['km_cluster_C_P_I']==2]['P_PRICE'])
7 plt.plot(train_data_salmon_cluster[train_data_salmon_cluster['km_cluster_C_P_I']==3]['P_PRICE'])
8
9
10 plt.legend(loc='best')
11 plt.show()

```

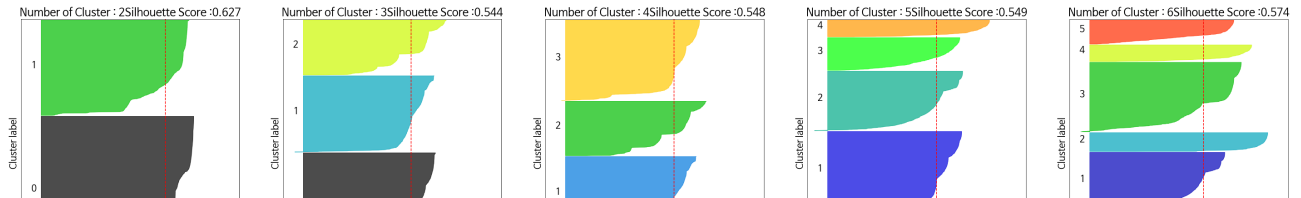
흰다리 새우 실루엣 검증

```

1 ## 실루엣 계수 사용을 위해 라벨인코딩
2 train_data_white_shrimp['sil_squid'] = le.fit_transform(train_data_white_shrimp['ctry_purpose_ir
3
4 ## 불필요한 열 제거
5 train_data_white_shrimp_sil = train_data_white_shrimp.drop(columns=['REG_DATE', 'CTRY_1', 'P_PUF

1 ## 실루엣 계수 시각화
2 ## 흰다리 새우 최적 클러스터 개수는 5개
3 ## 6개로 클러스터링 하면 실루엣계수는 높지만 0번 클러스터, 1번클러스터에서 평균값을 못넘는 데이
4
5 visualize_silhouette([2,3,4,5,6],train_data_white_shrimp_sil)

```



```

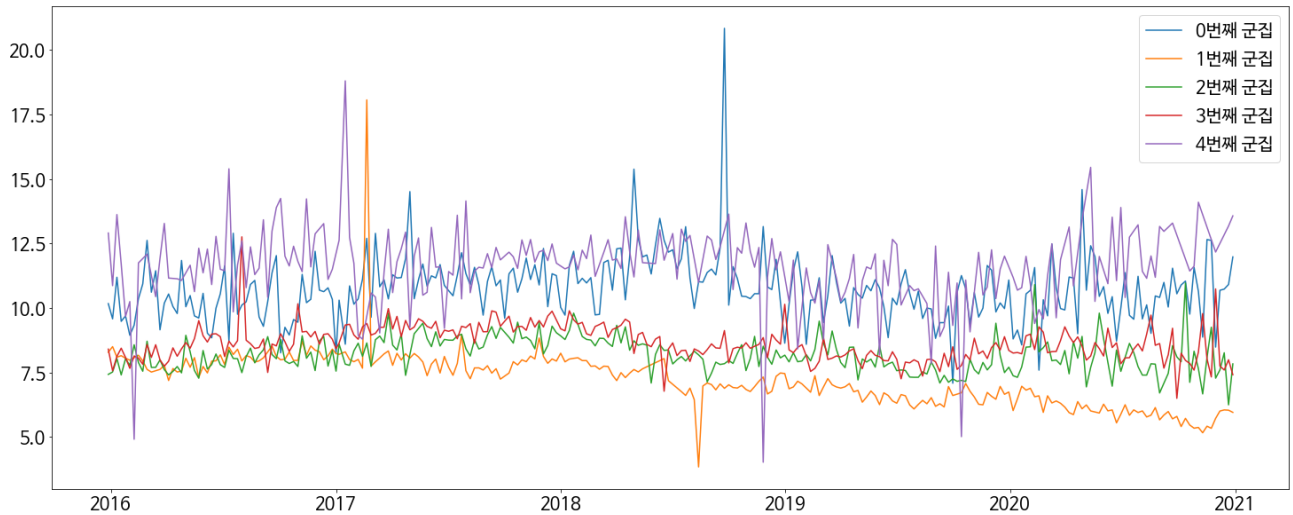
1 ## 군집화
2 def kmeans_data_white_shrimp(input_data):
3     df_C_P_I = input_data.groupby(['ctry_purpose_import_type', 'REG_DATE'])['P_PRICE'].mean().reset_index()
4     for i in range(len(df_C_P_I)):
5         df_C_P_I.iloc[i,1:] = (df_C_P_I.iloc[i,1:] - df_C_P_I.iloc[i,1:].mean())/df_C_P_I.iloc[i,1:].std()
6     df_C_P_I.fillna(0, inplace=True)
7     kmeans = KMeans(n_clusters=5, random_state = 7857)
8     km_cluster = kmeans.fit_predict(df_C_P_I.iloc[:,1:])
9
10    df_C_P_I_clust = df_C_P_I.copy()
11    df_C_P_I_clust['km_cluster_C_P_I'] = km_cluster
12
13    input_data=input_data.merge(df_C_P_I_clust[['ctry_purpose_import_type', 'km_cluster_C_P_I']])
14
15    return input_data
16

```

```

1 ## 흰다리새우 Kmeans 실행
2 train_data_white_shrimp=kmeans_data_white_shrimp(train_data_white_shrimp)
3
4 ## 흰다리 새우 시각화 준비
5 train_data_white_shrimp_cluster=train_data_white_shrimp.groupby(['REG_DATE', 'km_cluster_C_P_I']).mean().reset_index()
6 train_data_white_shrimp_cluster.index=train_data_white_shrimp_cluster['REG_DATE']
7
8
9
10 fig = plt.figure(figsize = (25, 10))
11
12 plt.plot(train_data_white_shrimp_cluster[train_data_white_shrimp_cluster['km_cluster_C_P_I']==0])
13 plt.plot(train_data_white_shrimp_cluster[train_data_white_shrimp_cluster['km_cluster_C_P_I']==1])
14 plt.plot(train_data_white_shrimp_cluster[train_data_white_shrimp_cluster['km_cluster_C_P_I']==2])
15 plt.plot(train_data_white_shrimp_cluster[train_data_white_shrimp_cluster['km_cluster_C_P_I']==3])
16 plt.plot(train_data_white_shrimp_cluster[train_data_white_shrimp_cluster['km_cluster_C_P_I']==4])
17
18
19 plt.legend(loc='best')
20 plt.show()

```



```

1 ## 상세어종별로 어떤 클러스터링 몇번씩 나타났는지 피벗 만들기위한 준비
2 train_data_squid_simple=train_data_squid[['REG_DATE', 'km_cluster_C_P_I']]
3 train_data_salmon_simple=train_data_salmon[['REG_DATE', 'km_cluster_C_P_I']]
4 train_data_white_shrimp_simple=train_data_white_shrimp[['REG_DATE', 'km_cluster_C_P_I']]
5
6
7 ## 상세어종별 피벗테이블
8 ## 오징어
9 train_data_squid_melt=train_data_squid_simple.melt(id_vars='REG_DATE')
10 train_data_squid_melt['var_name']=train_data_squid_melt.apply(lambda x : x['variable']+ '_' +str(x['km_cluster_C_P_I']), axis=1)
11 train_data_squid_pivot=pd.pivot_table(train_data_squid_melt, index='REG_DATE', columns='var_name')
12 columns=[i[2] for i in train_data_squid_pivot.columns]
13 train_data_squid_pivot.columns=columns
14 train_data_squid_pivot.reset_index(inplace=True)
15
16 ## 연어 train_data_salmon
17 train_data_salmon_melt=train_data_salmon_simple.melt(id_vars='REG_DATE')
18 train_data_salmon_melt['var_name']=train_data_salmon_melt.apply(lambda x : x['variable']+ '_' +str(x['km_cluster_C_P_I']), axis=1)
19 train_data_salmon_pivot=pd.pivot_table(train_data_salmon_melt, index='REG_DATE', columns='var_name')
20 columns=[i[2] for i in train_data_salmon_pivot.columns] ## 멀티 컬럼 제거
21 train_data_salmon_pivot.columns=columns
22 train_data_salmon_pivot.reset_index(inplace=True)
23
24 ##흰다리새우 train_data_white_shrimp
25 train_data_white_shrimp_melt=train_data_white_shrimp_simple.melt(id_vars='REG_DATE')
26 train_data_white_shrimp_melt['var_name']=train_data_white_shrimp_melt.apply(lambda x : x['variable']+ '_' +str(x['km_cluster_C_P_I']), axis=1)
27 train_data_white_shrimp_pivot=pd.pivot_table(train_data_white_shrimp_melt, index='REG_DATE', columns='var_name')
28 columns=[i[2] for i in train_data_white_shrimp_pivot.columns]
29 train_data_white_shrimp_pivot.columns=columns
30 train_data_white_shrimp_pivot.reset_index(inplace=True)

```

4. 2020년 데이터를 통한 모델 성능평가 및 2021년 클러스터링 예측

XGBoost 이용

```

1 ## 오징어 2015~2020년도 데이터로 2021년도 클러스터링 빈도 예측했을때 몇번씩 등장할것인지 예측
2 ## 학습데이터로 몇번째 전 주를 이용해서 학습시킬꺼지 확인
3 ## 13주전이 merror 가장 낮게 나옴 0.355
4
5 select_list_colname_squid=['km_cluster_C_P_I_0', 'km_cluster_C_P_I_1', 'km_cluster_C_P_I_2', 'km_
6
7 for input_before_week in range(5,14):
8     col_list=select_list_colname_squid
9     for i in col_list:
10         globals()['train_data_squid_pivot_{}'.format(i)] = train_data_squid_pivot[['REG_DATE',
11         for j in range(1,input_before_week):
12             globals()['train_data_squid_pivot_{}'.format(i)][str(j).zfill(2)+'주전']=train_data
13         globals()['train_data_squid_pivot_{}'.format(i)].dropna(inplace=True)
14
15     list_best_score=[]
16
17     ## 오징어
18     for i in tqdm(select_list_colname_squid):
19         data=globals()['train_data_squid_pivot_{}'.format(i)].copy()
20         data.reset_index(drop=True,inplace=True)
21
22         target=data[i]
23         data.drop('REG_DATE',axis=1,inplace=True)
24
25         len_data=data.shape[0]
26
27         feature_columns = list(data.columns.difference(select_list_colname_squid))
28
29         x = data[feature_columns]
30         y=target
31
32         x=x.applymap(lambda x : int(x))
33         y=y.apply(lambda x : int(x))
34
35         x.reset_index(drop=True,inplace=True)
36         y.reset_index(drop=True,inplace=True)
37
38         x_train, x_valid, y_train, y_valid = temporal_train_test_split(x, y, test_size = 0.2) #
39
40         dtrain=xgb.DMatrix(x_train,label=y_train)
41         dvalid=xgb.DMatrix(x_valid,label=y_valid)
42
43         params={'objective':'multi:softmax','eval_metric':'merror','random_state':7857,'num_clas
44
45         num_round=500
46         watchlist=[(dtrain,'train'),(dvalid,'eval')]
47         model=xgb.train(params,dtrain,num_round,evals=watchlist,early_stopping_rounds=20,verbose
48
49         for j in range(0,27):
50             dtest=xgb.DMatrix(pd.DataFrame(data.iloc[-1,:].values[0:-1].reshape(1,-1),columns=fe
51             result=pd.DataFrame(data.iloc[-1,:].values[0:-1].reshape(1,-1),columns=feature_colur

```

```

52     result.insert(0,i,int(np.round(model.predict(dtest),0)))
53     data=data.append(result)
54     data.reset_index(drop=True,inplace=True)
55
56
57     globals()['list_cl_{}'.format(i)]=list(data[i].iloc[len_data:].values)
58
59     list_best_score.append(model.best_score)
60
61     print(input_before_week,'주전 변수 있을때 평균 eval-merror:',np.mean(list_best_score))

```

```

100% 4/4 [00:00<00:00, 4.88it/s]

```

```

5 주전 변수 있을때 평균 eval-merror: 0.35576925

```

```

100% 4/4 [00:00<00:00, 5.24it/s]

```

```

6 주전 변수 있을때 평균 eval-merror: 0.35784324999999995

```

```

100% 4/4 [00:00<00:00, 5.21it/s]

```

```

7 주전 변수 있을때 평균 eval-merror: 0.382353

```

```

100% 4/4 [00:00<00:00, 4.66it/s]

```

```

8 주전 변수 있을때 평균 eval-merror: 0.377451

```

```

100% 4/4 [00:00<00:00, 4.53it/s]

```

```

9 주전 변수 있을때 평균 eval-merror: 0.39705874999999996

```

```

100% 4/4 [00:00<00:00, 4.64it/s]

```

```

10 주전 변수 있을때 평균 eval-merror: 0.36764725

```

```

100% 4/4 [00:00<00:00, 4.50it/s]

```

```

11 주전 변수 있을때 평균 eval-merror: 0.40499999999999997

```

```

100% 4/4 [00:00<00:00, 4.53it/s]

```

```

12 주전 변수 있을때 평균 eval-merror: 0.415

```

```

100% 4/4 [00:00<00:00, 4.58it/s]

```

```

13 주전 변수 있을때 평균 eval-merror: 0.355

```

```

1 ## 2021년 오징어 클러스터링 예측

```

```

2 future_squid_cluster=pd.DataFrame({'km_cluster_C_P_I_0':list_cl_km_cluster_C_P_I_0,
3                                     'km_cluster_C_P_I_1':list_cl_km_cluster_C_P_I_1,
4                                     'km_cluster_C_P_I_2':list_cl_km_cluster_C_P_I_2,
5                                     'km_cluster_C_P_I_3':list_cl_km_cluster_C_P_I_3})

```

```

1 ## 연어 15주전 0.145

```

```

2 select_list_colname_salmon=['km_cluster_C_P_I_0', 'km_cluster_C_P_I_1', 'km_cluster_C_P_I_2', 'kr
3
4
5 # 연어
6
7 for input_before_week in range(5,16):
8     col_list=select_list_colname_salmon
9     for i in col_list:
10         globals()['train_data_salmon_pivot_{}'.format(i)] = train_data_salmon_pivot[['REG_DATE
11         for j in range(1,input_before_week):
12             globals()['train_data_salmon_pivot_{}'.format(i)][str(j).zfill(2)+'주전']=train_da

```

```

13     globals()['train_data_salmon_pivot_{}'.format(i)].dropna(inplace=True)
14     list_best_score=[]
15
16     ## 연어
17     for i in tqdm(select_list_colname_salmon):
18         data=globals()['train_data_salmon_pivot_{}'.format(i)].copy()
19         data.reset_index(drop=True,inplace=True)
20
21         target=data[i]
22         data.drop('REG_DATE',axis=1,inplace=True)
23
24         len_data=data.shape[0]
25
26         feature_columns = list(data.columns.difference(select_list_colname_salmon))
27
28         x = data[feature_columns]
29         y=target
30
31         x=x.applymap(lambda x : int(x))
32         y=y.apply(lambda x : int(x))
33
34         x.reset_index(drop=True,inplace=True)
35         y.reset_index(drop=True,inplace=True)
36
37         x_train, x_valid, y_train, y_valid = temporal_train_test_split(x, y, test_size = 0.2)
38
39         dtrain=xgb.DMatrix(x_train,label=y_train)
40         dvalid=xgb.DMatrix(x_valid,label=y_valid)
41
42         params={'objective':'multi:softmax','eval_metric':'merror','random_state':7857,'num_class':4}
43
44         num_round=500
45         watchlist=[(dtrain,'train'),(dvalid,'eval')]
46         model=xgb.train(params,dtrain,num_round,evals=watchlist,early_stopping_rounds=20,verbose=False)
47
48         for j in range(0,27):
49             dtest=xgb.DMatrix(pd.DataFrame(data.iloc[-1,:].values[0:-1].reshape(1,-1),columns=feature_columns))
50             result=pd.DataFrame(data.iloc[-1,:].values[0:-1].reshape(1,-1),columns=feature_columns)
51             result.insert(0,i,int(np.round(model.predict(dtest),0)))
52             data=data.append(result)
53             data.reset_index(drop=True,inplace=True)
54
55
56         globals()['list_cl_{}'.format(i)]=list(data[i].iloc[len_data:].values)
57
58         list_best_score.append(model.best_score)
59
60     print(input_before_week,'주전 변수 있을때 평균 eval-merror:',np.mean(list_best_score))

```

```
25 feature_columns = list(data.columns.difference(select_list_colname_white_shrimp))
```

```

26
27     x = data[feature_columns]
28     y=target
29
30     x=x.applymap(lambda x : int(x))
31     y=y.apply(lambda x : int(x))
32
33     x.reset_index(drop=True,inplace=True)
34     y.reset_index(drop=True,inplace=True)
35
36
37     x_train, x_valid, y_train, y_valid = temporal_train_test_split(x, y, test_size = 0.2)
38
39     dtrain=xgb.DMatrix(x_train,label=y_train)
40     dvalid=xgb.DMatrix(x_valid,label=y_valid)
41
42     params={'objective':'multi:softmax','eval_metric':'merror','random_state':7857,'num_class':4}
43
44     num_round=500
45     watchlist=[(dtrain,'train'),(dvalid,'eval')]
46     model=xgb.train(params,dtrain,num_round,evals=watchlist,early_stopping_rounds=20,verbose=0)
47
48
49     for j in range(0,27):
50         dtest=xgb.DMatrix(pd.DataFrame(data.iloc[-1,:].values[0:-1].reshape(1,-1),columns=feature_columns))
51         result=pd.DataFrame(data.iloc[-1,:].values[0:-1].reshape(1,-1),columns=feature_columns)
52         # result.insert(0,'km_cluster_C_P_I_{}'.format(i),int(np.round(model.predict(dtest))))
53         result.insert(0,i,int(np.round(model.predict(dtest),0)))
54         data=data.append(result)
55         data.reset_index(drop=True,inplace=True)
56
57
58     globals()['list_cl_{}'.format(i)]=list(data[i].iloc[len_data:].values)
59
60     list_best_score.append(model.best_score)
61
62     print(input_before_week,'주전 변수 있을때 평균 eval-merror:',np.mean(list_best_score))

```


100% 5/5 [00:00<00:00, 5.73it/s]

5 주전 변수 있을때 평균 eval-merror: 0.21923079999999998

100% 5/5 [00:00<00:00, 5.76it/s]

6 주전 변수 있을때 평균 eval-merror: 0.22692299999999999

100% 5/5 [00:00<00:00, 5.62it/s]

7 주전 변수 있을때 평균 eval-merror: 0.22307600000000000

```
1 future_white_shrimp_cluster=pd.DataFrame({'km_cluster_C_P_I_0':list_cl_km_cluster_C_P_I_0,
2                                           'km_cluster_C_P_I_1':list_cl_km_cluster_C_P_I_1,
3                                           'km_cluster_C_P_I_2':list_cl_km_cluster_C_P_I_2,
4                                           'km_cluster_C_P_I_3':list_cl_km_cluster_C_P_I_3,
5                                           'km_cluster_C_P_I_4':list_cl_km_cluster_C_P_I_4})
```

100% 5/5 [00:01<00:00, 5.25it/s]

▼ 5. 외부데이터 전처리

11 주전 변수 있을때 평균 eval-merror: 0.2200290

▼ 환율데이터

100% 5/5 [00:01<00:00, 4.07it/s]

```
1 ## 환율 데이터 경로
2 ex_can = './외부데이터/USD_CAD 내역.csv' ## 캐나다
3 ex_CHL = './외부데이터/USD_CLP 내역.csv' ## 칠레
4 ex_CH = './외부데이터/USD_CNH 내역.csv' ## 중국
5 ex_ENG = './외부데이터/USD_GBP 내역.csv' ## 영국
6 ex_KOR = './외부데이터/USD_KRW 내역.csv' ## 대한민국
7 ex_MAL = './외부데이터/USD_MYR 내역.csv' ## 말레이시아
8 ex_NOR = './외부데이터/USD_NOK 내역.csv' ## 노르웨이
9 ex_PE = './외부데이터/USD_PEN 내역.csv' ## 페루
10 ex_TH = './외부데이터/USD_THB 내역.csv' ## 태국
11 ex_VN = './외부데이터/USD_VND 내역.csv' ## 베트남
12
13
14 ## 국가별 환율 데이터 불러오기
15 train_full_data=pd.DataFrame()
16 ex_CAN=pd.read_csv(ex_can, parse_dates=['날짜'])
17 ex_CHL=pd.read_csv(ex_CHL, parse_dates=['날짜'])
18 ex_CH=pd.read_csv(ex_CH, parse_dates=['날짜'])
19 ex_ENG=pd.read_csv(ex_ENG, parse_dates=['날짜'])
20 ex_KOR=pd.read_csv(ex_KOR, parse_dates=['날짜'])
21 ex_MAL=pd.read_csv(ex_MAL, parse_dates=['날짜'])
22 ex_NOR=pd.read_csv(ex_NOR, parse_dates=['날짜'])
23 ex_PE=pd.read_csv(ex_PE, parse_dates=['날짜'])
24 ex_TH=pd.read_csv(ex_TH, parse_dates=['날짜'])
25 ex_VN=pd.read_csv(ex_VN, parse_dates=['날짜'])
```

```
1 # 문자열 형태의 날짜 칼럼 datetime 형식으로 변환
2 def str_to_date(str_1):
3
4     str_2 = str_1[:4] + '-' + str_1[6:8]+'-'+str_1[10:12] # 'YYYY년 MM월 DD일' 형태의 데이터를
5
```

```

6     date_return = datetime.strptime(str_2, '%Y-%m-%d')
7     return date_return
8
9 # 문자열 형태의 증가 칼럼을 분석에 맞게 float형태로 변환
10 # 대한민국용
11 def str_to_float_kor(str_1):
12
13     str_2 = str_1[0] +str_1[2:] # 0,000.00 형태에서 1000의자리 소수점 제거
14
15     rate_return = float(str_2)
16     return rate_return
17
18 # 베트남용
19 def str_to_float_vn(str_1):
20
21     str_2 = str_1[:2] +str_1[3:] # 00,000.00 형태에서 1000의자리 소수점 제거
22
23     rate_return = float(str_2)
24     return rate_return


1 ## 국가 데이터프레임 이름으로 된 리스트 선언
2 ctry = ['ex_CAN', 'ex_CHL', 'ex_CH', 'ex_ENG', 'ex_KOR', 'ex_MAL', 'ex_NOR', 'ex_PE', 'ex_TH', 'ex_VN']
3
4 ## 각 국가 데이터프레임의 날짜칼럼 datetime 데이터타입으로 변환
5 for i in ctry:
6     globals()[i]['날짜'] = globals()[i]['날짜'].apply(lambda x: str_to_date(x))
7
8 ## 증가 칼럼이 문자열로 이루어진 한국과 베트남 데이터프레임에 대해서 문자열 데이터를 float형 데이터로 변환
9 ex_KOR['증가'] = ex_KOR['증가'].apply(lambda x: str_to_float_kor(x))
10 ex_VN['증가'] = ex_VN['증가'].apply(lambda x: str_to_float_vn(x))
11
12 ## 대한민국데이터의 증가 칼럼을 KRW증가로 변경
13 ex_KOR.rename(columns={'증가': 'KRW증가'}, inplace=True)
14
15 ## 각 데이터에서 불필요한 칼럼(오픈, 고가, 저가, 변동) 제거
16 for i in ctry:
17     globals()[i].drop(columns=['오픈', '고가', '저가', '변동 %'], inplace=True)
18
19
20 # 모델에서 분석할 대한민국을 제외한 나머지 국가들 데이터만 추출
21 ctry_concat = ['ex_CAN', 'ex_CHL', 'ex_CH', 'ex_ENG', 'ex_MAL', 'ex_NOR', 'ex_PE', 'ex_TH', 'ex_VN']
22
23 ## 기존의 달러대비 환율을 통해 원화대비 환율을 새로 계산한 후 새로운 칼럼으로 생성
24 for i in ctry_concat:
25     globals()[i] = pd.merge(globals()[i], ex_KOR, how='left', left_on='날짜', right_on='날짜')
26     globals()[i]['원화대비 환율'+ '_' +i[3:]] = globals()[i]['KRW증가']/globals()[i]['증가']

```

▼ 유가 데이터

```

1 ## 유가 데이터 경로 및 불러오기
2 oil = './외부데이터/국제유가2021-09-09.csv' ## 유가
3 oil_price = pd.read_csv(oil, parse_dates=['Date'])

```

▼ 외부 데이터 전처리 작업

```

1 ## 각 어종별 데이터프레임을 따로 생성
2
3 ## 오징어
4 ## 페루, 칠레, 중국
5 ctry_squid=['ex_PE','ex_CHL','ex_CH']
6 for i in ctry_squid:
7     globals()[i].drop(columns=['종가','KRW종가'],inplace=True)          # 불필요한 칼럼 제거
8
9 ## 연어
10 ## 노르웨이, 영국, 캐나다
11 ctry_salmon=['ex_NOR','ex_ENG','ex_CAN']
12 for i in ctry_salmon:
13     globals()[i].drop(columns=['종가','KRW종가'],inplace=True)          # 불필요한 칼럼 제거
14
15 ## 흰다리새우
16 ## 베트남, 태국, 말레이시아
17 ctry_white_shrimp=['ex_VN','ex_TH','ex_MAL']
18 for i in ctry_white_shrimp:
19     globals()[i].drop(columns=['종가','KRW종가'],inplace=True)          # 불필요한 칼럼 제거
20

```

```

1 ## 어종에 따라 같은 어종 국가는 하나로 병합
2
3 ## 오징어 주요 수출국 환율데이터
4 ex_squid = pd.merge(ex_PE,ex_CHL, how='left',left_on='날짜', right_on='날짜')
5 ex_squid = pd.merge(ex_squid,ex_CH, how='left',left_on='날짜', right_on='날짜')
6 ex_squid = pd.merge(ex_squid,ex_KOR, how='left',left_on='날짜', right_on='날짜')
7
8 ## 연어 주요 수출국 환율데이터
9 ex_salmon = pd.merge(ex_NOR,ex_ENG, how='left',left_on='날짜', right_on='날짜')
10 ex_salmon = pd.merge(ex_salmon,ex_CAN, how='left',left_on='날짜', right_on='날짜')
11 ex_salmon = pd.merge(ex_salmon,ex_KOR, how='left',left_on='날짜', right_on='날짜')
12
13
14 ## 흰다리새우 주요 수출국 환율데이터
15 ex_white_shrimp = pd.merge(ex_VN,ex_TH, how='left',left_on='날짜', right_on='날짜')
16 ex_white_shrimp = pd.merge(ex_white_shrimp,ex_MAL, how='left',left_on='날짜', right_on='날짜')
17 ex_white_shrimp = pd.merge(ex_white_shrimp,ex_KOR, how='left',left_on='날짜', right_on='날짜')
18
19
20 ## 오징어
21 ## 오징어 피벗 테이블과 유가데이터 병합
22 data_squid_oil = pd.merge(train_data_squid_pivot, oil_price, how='left', left_on='REG_DATE', right_on='REG_DATE')
23
24 ## 오징어 피벗 테이블 + 유가 + 환율(오징어 주요 수출국)
25 data_squid_oil_ex = pd.merge(data_squid_oil, ex_squid, how='left', left_on='REG_DATE', right_on='날짜')
26
27
28

```

```

9 ## 연어
10 ## 연어 피벗 테이블과 유가데이터 병합
11 data_salmon_oil = pd.merge(train_data_salmon_pivot, oil_price, how='left', left_on='REG_DATE', right_on='REG_DATE')
12
13 ## 연어 피벗 테이블 + 유가 + 환율(연어 주요 수출국)
14 data_salmon_oil_ex = pd.merge(data_salmon_oil, ex_salmon, how='left', left_on='REG_DATE', right_on='REG_DATE')
15
16
17 ## 흰다리새우
18 ## 흰다리새우 피벗 테이블과 유가데이터 병합
19 data_white_shrimp_oil = pd.merge(train_data_white_shrimp_pivot, oil_price, how='left', left_on='REG_DATE', right_on='REG_DATE')
20
21 ## 흰다리새우 피벗 테이블 + 유가 + 환율(연어 주요 수출국)
22 data_white_shrimp_oil_ex = pd.merge(data_white_shrimp_oil, ex_white_shrimp, how='left', left_on='REG_DATE', right_on='REG_DATE')

1 ## 각 어종별 데이터에서 REG_DATE를 제외한 datetime 타입 칼럼 제거
2
3 data_squid_oil_ex.drop(columns=['Date', '날짜'], inplace=True)
4 data_salmon_oil_ex.drop(columns=['Date', '날짜'], inplace=True)
5 data_white_shrimp_oil_ex.drop(columns=['Date', '날짜'], inplace=True)

1 ## 결측치 처리
2 data_squid_oil_ex.interpolate(method="pad", inplace=True)
3 data_salmon_oil_ex.interpolate(method="pad", inplace=True)
4 data_white_shrimp_oil_ex.interpolate(method="pad", inplace=True)

```

6. 제품별 평균 단가 예측 성능 평가

target data : 2020년도 데이터 성능파악결과 예측 방법론에 문제 없다고 판단하여 2021년도 데이터 예측 진행

```

1 ## 가격예측에 맞게 불필요한 열 제거
2 train_data_squid_cluster_price_grp = train_data_squid.drop(columns=['CTRY_1', 'P_PURPOSE', 'P_NAME'])
3 train_data_salmon_cluster_price_grp = train_data_salmon.drop(columns=['CTRY_1', 'P_PURPOSE', 'P_NAME'])
4 train_data_white_shrimp_cluster_price_grp = train_data_white_shrimp.drop(columns=['CTRY_1', 'P_PURPOSE', 'P_NAME'])

1 ## 평가용 데이터 날짜 불러오기
2 REG_DATE_2021=pd.read_excel('./제공데이터/2021 빅콘테스트_데이터분석분야_챔피언리그_수산Biz_평가용_데이터.xlsx')
3 eval=REG_DATE_2021[['일자']]
4 eval.rename(columns={'일자':'REG_DATE'}, inplace=True)

1 ## 2021년 클러스터링
2 eval_squid = pd.merge(eval, future_squid_cluster, how='left', left_index=True, right_index=True)
3 eval_salmon = pd.merge(eval, future_salmon_cluster, how='left', left_index=True, right_index=True)
4 eval_white_shrimp = pd.merge(eval, future_white_shrimp_cluster, how='left', left_index=True, right_index=True)

1 ## 유가 및 환율 데이터 연결
2 eval_squid_oil = pd.merge(eval_squid, oil_price, how='left', left_on='REG_DATE', right_on='REG_DATE')
3 eval_salmon_oil = pd.merge(eval_salmon, oil_price, how='left', left_on='REG_DATE', right_on='REG_DATE')
4 eval_white_shrimp_oil = pd.merge(eval_white_shrimp, oil_price, how='left', left_on='REG_DATE', right_on='REG_DATE')

```

```

3 eval_squid_oil_ex = pd.merge(eval_squid_oil, ex_squid, how='left', left_on='REG_DATE', right_on='
4
5 eval_salmon_oil = pd.merge(eval_salmon, oil_price, how='left', left_on='REG_DATE', right_on='Date
6 eval_salmon_oil_ex = pd.merge(eval_salmon_oil, ex_salmon, how='left', left_on='REG_DATE', right
7
8 eval_white_shrimp_oil = pd.merge(eval_white_shrimp, oil_price, how='left', left_on='REG_DATE', r
9 eval_white_shrimp_oil_ex = pd.merge(eval_white_shrimp_oil, ex_white_shrimp, how='left', left_on
10
11 eval_squid_oil_ex.drop(columns=['Date', '날짜'], inplace=True)
12 eval_salmon_oil_ex.drop(columns=['Date', '날짜'], inplace=True)
13 eval_white_shrimp_oil_ex.drop(columns=['Date', '날짜'], inplace=True)

```

```

1 ## 2015 ~ 2020년도 데이터에 가격 붙이기
2 train_data_squid_cluster_price_grp = train_data_squid.drop(columns=['CTRY_1', 'P_PURPOSE', 'P_NA
3 train_data_salmon_cluster_price_grp = train_data_salmon.drop(columns=['CTRY_1', 'P_PURPOSE', 'P_I
4 train_data_white_shrimp_cluster_price_grp = train_data_white_shrimp.drop(columns=['CTRY_1', 'P_P
5
6 data_squid_oil_ex_price = pd.merge(data_squid_oil_ex, train_data_squid_cluster_price_grp, how='
7 data_salmon_oil_ex_price = pd.merge(data_salmon_oil_ex, train_data_salmon_cluster_price_grp, ho
8 data_white_shrimp_oil_ex_price = pd.merge(data_white_shrimp_oil_ex, train_data_white_shrimp_clus
9

```

```

1 ## 결측치 보간법으로 대체
2 data_squid_oil_ex_price.interpolate(method="pad", inplace=True)
3 data_salmon_oil_ex_price.interpolate(method="pad", inplace=True)
4 data_white_shrimp_oil_ex_price.interpolate(method="pad", inplace=True)

```

▼ 가격 예측 모델링

```

1 ## mse 계산하는 생성함수
2 def get_rmse(y, pred):
3     mse = mean_squared_error(y, pred)
4     rmse = np.sqrt(mse)
5     print('RMSE 값:', np.round(rmse, 3))
6     return rmse

1 ## 2021년도까지 연결한 데이터에서 2021년도 제거
2 data_squid_price = pd.concat([data_squid_oil_ex_price, eval_squid_oil_ex], ignore_index=True)
3 data_salmon_price = pd.concat([data_salmon_oil_ex_price, eval_salmon_oil_ex], ignore_index=True)
4 data_white_shrimp_price = pd.concat([data_white_shrimp_oil_ex_price, eval_white_shrimp_oil_ex],
5
6 data_squid_price_test=data_squid_price.copy()
7 data_salmon_price_test=data_salmon_price.copy()
8 data_white_shrimp_price_test=data_white_shrimp_price.copy()
9
10 data_squid_price.dropna(axis=0, inplace=True)
11 data_salmon_price.dropna(axis=0, inplace=True)
12 data_white_shrimp_price.dropna(axis=0, inplace=True)

```

오징어 가격 예측 성능평가

```

1 ## 1주 ~ 12주 전 가격을 각각 하나의 칼럼을 생성하여 저장
2 for i in range(1,13):
3     data_squid_price['{}주전 금액'.format(i)]=data_squid_price['P_PRICE'].shift(i)
4
5
6 rmse_list =[] ## 평균 rmse 계산
7
8
9 ## 2020년 이전 데이터는 Train 데이터로
10 ## 2020년 데이터는 Test 데이터로 저장
11 data_squid_20_train = data_squid_price.head(-52)
12 data_squid_20_test = data_squid_price.tail(52)
13
14 ## 2020년 가격 예측 결과를 저장할 P_PRICE_pred
15 data_squid_20_train['P_PRICE_pred'] = data_squid_20_train['P_PRICE']
16 data_squid_20_test['P_PRICE_pred'] = np.nan # 2020년 데이터의 P_PRICE_pred 칼럼은 모두
17
18 # 2020년 데이터를 한 행씩 예측 후 Train 데이터에 연결 후 연결된 새로운 데이터를 바탕으로 다시 예
19 for i in tqdm(range(208,260)):
20     test = data_squid_20_test[data_squid_20_test.index == i]
21
22
23 X_train = data_squid_20_train.drop(['P_PRICE', 'REG_DATE', 'P_PRICE_pred'], axis=1, inplace=False)
24 y_train = data_squid_20_train['P_PRICE_pred']
25
26 X_test = data_squid_20_test.drop(['P_PRICE', 'REG_DATE', 'P_PRICE_pred'], axis=1, inplace=False)
27 y_test = data_squid_20_test['P_PRICE']
28
29 my_model = XGBRegressor(n_estimators=10000, learning_rate=0.01, seed=7857, max_depth=30, colsample_bytree=0.8)
30 my_model.fit(X_train, y_train, verbose=False)
31 pred = my_model.predict(X_test)
32 data_squid_20_test['P_PRICE_pred'] = pred # 예측 결과 값은 P_PRICE_pred 칼럼에 저장
33
34 data_squid_20_train = pd.concat([data_squid_20_train, data_squid_20_test])
35 rmse=get_rmse(y_test, pred)
36 rmse_list.append(rmse) # 각 행마다 RMSE 출력
37
38 for i in range(1,13):
39     data_squid_20_train['{}주전 금액'.format(i)]=data_squid_20_train['P_PRICE_pred'].shift(i)
40
41
42
43 rmse_avg = np.array(rmse_list)
44 print('평균 RMSE 값은 :',rmse_avg.mean())

```

100%

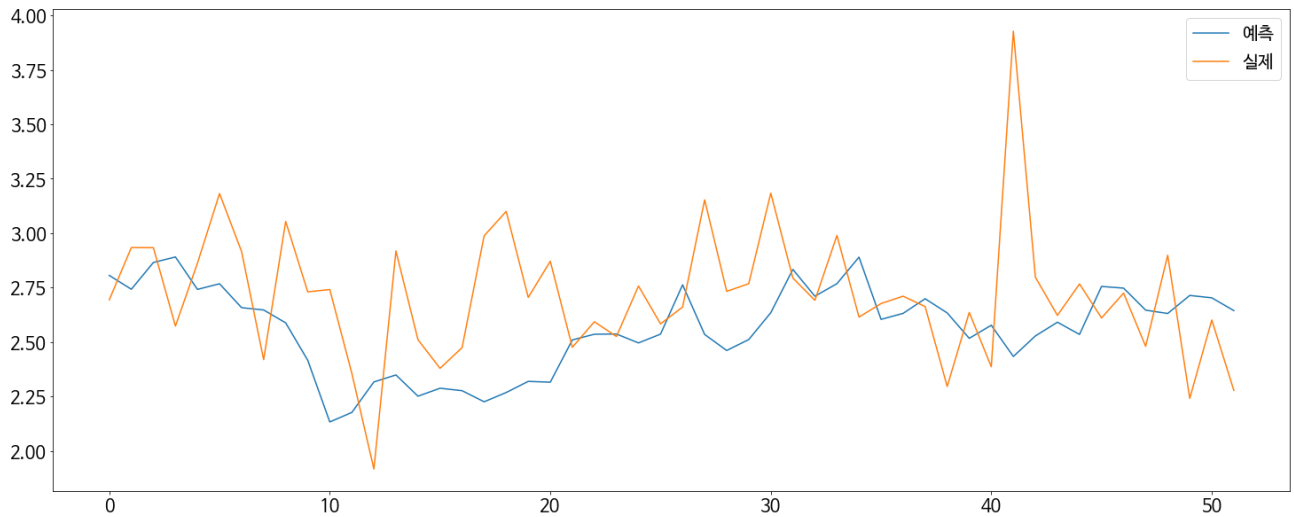
52/52 [08:31<00:00, 17.73s/it]

RMSE 값: 0.433
 RMSE 값: 0.402
 RMSE 값: 0.395
 RMSE 값: 0.395
 RMSE 값: 0.389
 RMSE 값: 0.392
 RMSE 값: 0.392
 RMSE 값: 0.388
 RMSE 값: 0.387
 RMSE 값: 0.387
 RMSE 값: 0.387
 RMSE 값: 0.385
 RMSE 값: 0.382
 RMSE 값: 0.381
 RMSE 값: 0.383
 RMSE 값: 0.385
 RMSE 값: 0.387
 RMSE 값: 0.389
 RMSE 값: 0.387
 RMSE 값: 0.39
 RMSE 값: 0.389
 RMSE 값: 0.39
 RMSE 값: 0.39
 RMSE 값: 0.389
 RMSE 값: 0.389
 RMSE 값: 0.39
 RMSE 값: 0.392
 RMSE 값: 0.392
 RMSE 값: 0.39
 RMSE 값: 0.389
 RMSE 값: 0.39
 RMSE 값: 0.388
 RMSE 값: 0.387
 RMSE 값: 0.388
 RMSE 값: 0.39
 RMSE 값: 0.387
 RMSE 값: 0.388
 RMSE 값: 0.387
 RMSE 값: 0.386
 RMSE 값: 0.386
 RMSE 값: 0.387
 RMSE 값: 0.386
 RMSE 값: 0.385
 RMSE 값: 0.385
 RMSE 값: 0.385
 RMSE 값: 0.384

```

1 ## 시각화
2 fig = plt.figure(figsize = (25, 10))
3
4 plt.plot(my_model.predict(X_test), label='예측')
5 plt.plot(y_test.values, label='실제')
6
7 plt.legend(loc='best')
8 plt.show()
9

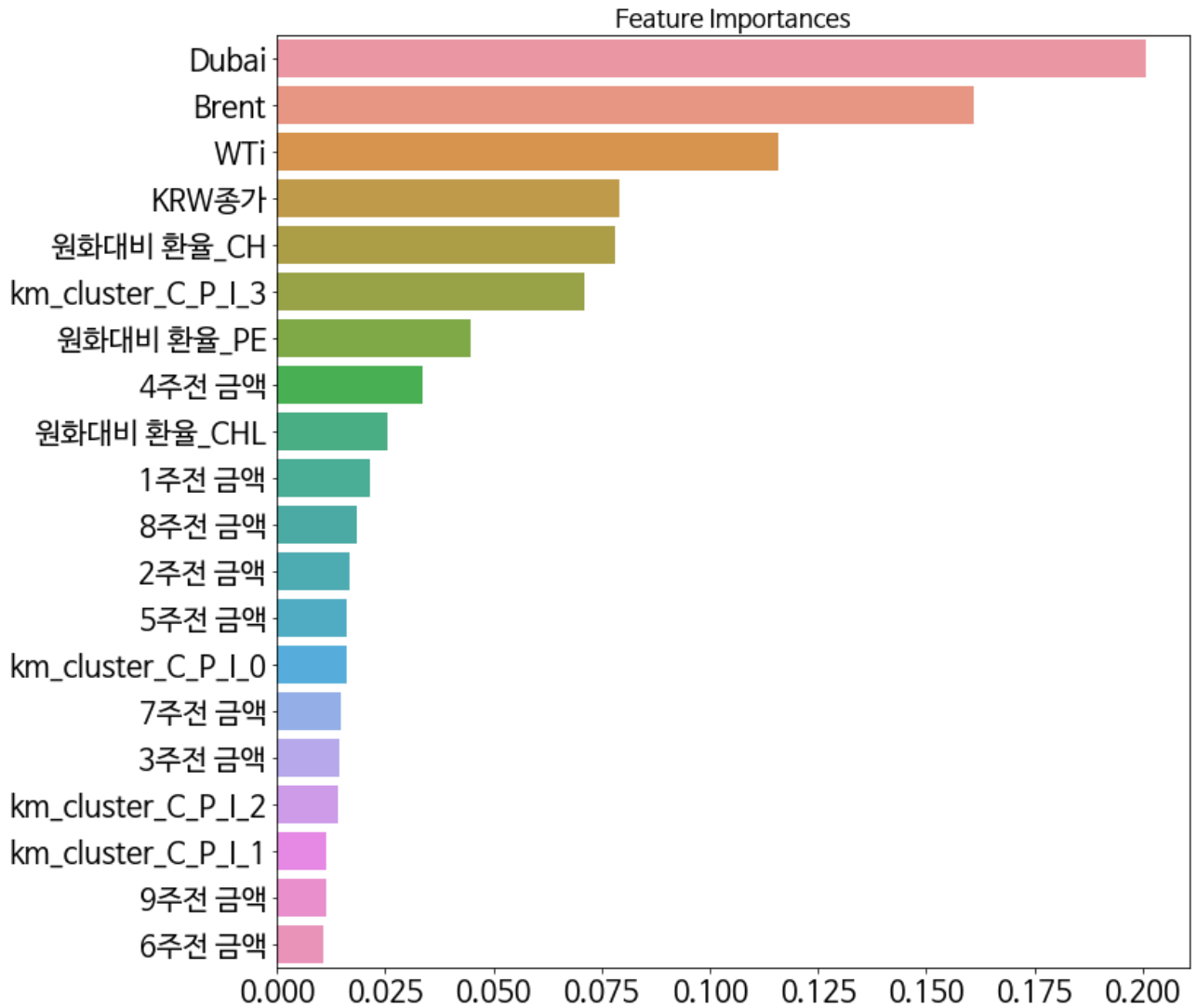
```



```

1 ## 중요도 파악
2 ## 외부데이터가 중요한것으로 판단
3 def get_top_features(model):
4     ftr_importances_values = model.feature_importances_
5     ftr_importances = pd.Series(ftr_importances_values, index=X_train.columns )
6     ftr_top20 = ftr_importances.sort_values(ascending=False)[:20]
7     return ftr_top20
8
9 def visualize_ftr_importances(models):
10     # 2개 회귀 모델의 시각화를 위해 2개의 컬럼을 가지는 subplot 생성
11     fig, axs = plt.subplots(figsize=(10,10),nrows=1, ncols=1)
12     fig.tight_layout()
13     # 입력인자로 받은 list객체인 models에서 차례로 model을 추출하여 피쳐 중요도 시각화.
14
15     # 중요도 상위 20개의 피쳐명과 그때의 중요도값 추출
16     ftr_top20 = get_top_features(models)
17     plt.title('Feature Importances',size=17)
18     sns.barplot(x=ftr_top20.values, y=ftr_top20.index ) # , ax=axs[1]
19
20
21 visualize_ftr_importances(my_model)

```

연어 가격 예측 성능평가

```

1
2 ## 1주 ~ 12주 전 가격을 각각 하나의 칼럼을 생성하여 저장
3 for i in range(1,13):
4     data_salmon_price['{}주전 금액'.format(i)]=data_salmon_price['P_PRICE'].shift(i)
5
6 rmse_list =[]
7
8 ## 2020년 이전 데이터는 Train 데이터로
9 ## 2020년 데이터는 Test 데이터로 저장
10 data_salmon_20_train = data_salmon_price.head(-50)
11 data_salmon_20_test = data_salmon_price.tail(50)
12
13 ## 2020년 가격 예측 결과를 저장할 P_PRICE_pred
14 data_salmon_20_train['P_PRICE_pred'] = data_salmon_20_train['P_PRICE']
15 data_salmon_20_test['P_PRICE_pred'] = np.nan          # 2020년 데이터의 P_PRICE_pred 칼럼은 모두
16

```

```

17 # 2020년 데이터를 한 행씩 예측 후 Train 데이터에 연결 후 연결된 새로운 데이터를 바탕으로 다시 여
18 for i in tqdm(range(210,260)):
19     test = data_salmon_20_test[data_salmon_20_test.index == i]
20
21
22     X_train = data_salmon_20_train.drop(['P_PRICE', 'REG_DATE', 'P_PRICE_pred'], axis=1, inplace=False)
23     y_train = data_salmon_20_train['P_PRICE']
24
25     X_test = data_salmon_20_test.drop(['P_PRICE', 'REG_DATE', 'P_PRICE_pred'], axis=1, inplace=False)
26     y_test = data_salmon_20_test['P_PRICE']
27
28     my_model = XGBRegressor(n_estimators=10000, learning_rate=0.01, seed=7857, max_depth=30, colsample_bytree=0.9)
29     my_model.fit(X_train, y_train, verbose=False)
30     pred = my_model.predict(X_test)
31     data_salmon_20_test['P_PRICE_pred'] = pred # 예측 결과 값은 P_PRICE_pred 칼럼에 저장
32
33     data_salmon_20_train = pd.concat([data_salmon_20_train, data_salmon_20_test])
34     rmse=get_rmse(y_test, pred)
35     rmse_list.append(rmse) # 각 행마다 RMSE 출력
36
37     for i in range(1,13):
38         data_salmon_20_train['{}주전 금액'.format(i)]=data_salmon_20_train['P_PRICE_pred'].shift(i)
39
40 rmse_avg = np.array(rmse_list)
41 print('평균 RMSE 값은 :',rmse_avg.mean())

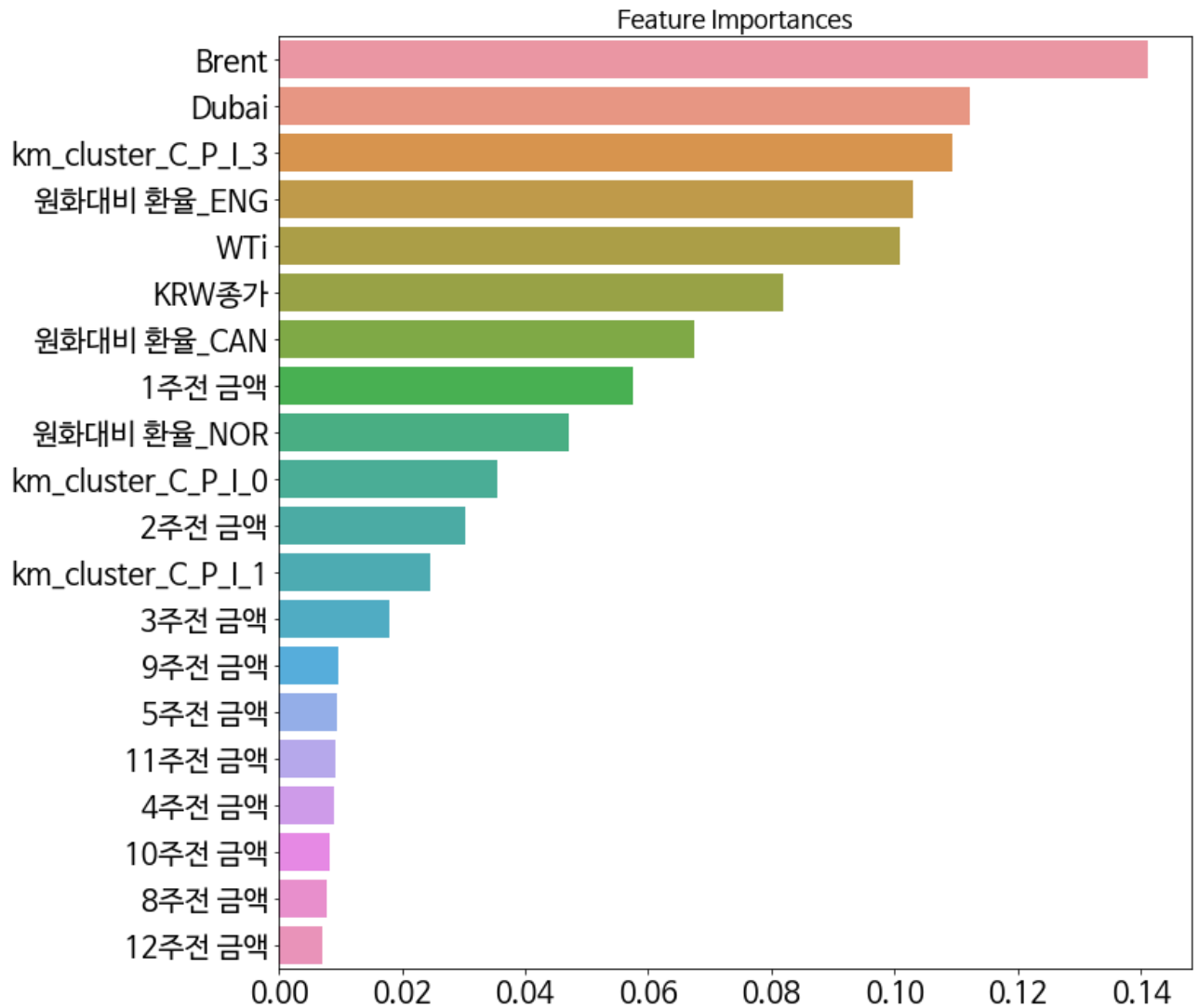
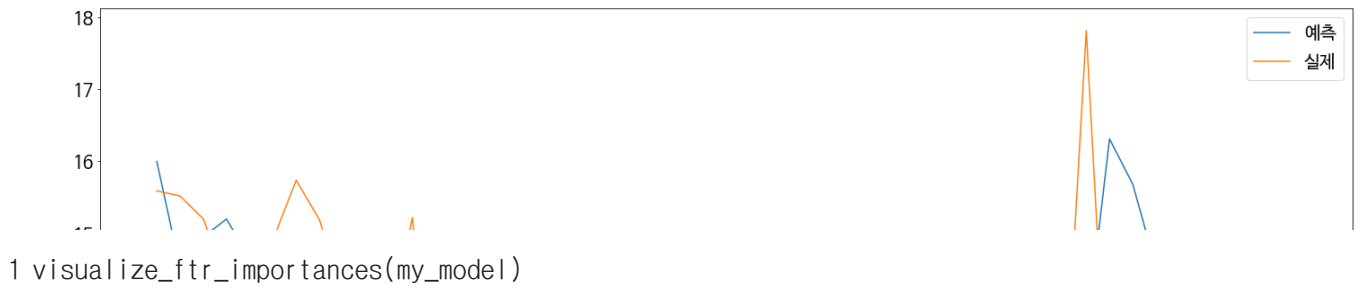
```

100%

50/50 [08:37<00:00, 18.54s/it]

RMSE 값: 1.14
RMSE 값: 1.196
RMSE 값: 1.214
RMSE 값: 1.236
RMSE 값: 1.249
RMSE 값: 1.246
RMSE 값: 1.241
RMSE 값: 1.256
RMSE 값: 1.246
RMSE 값: 1.261
RMSE 값: 1.245
RMSE 값: 1.251
RMSE 값: 1.245
RMSE 값: 1.254
RMSE 값: 1.26
RMSE 값: 1.274
RMSE 값: 1.27
RMSE 값: 1.269
RMSE 값: 1.268
RMSE 값: 1.276
RMSE 값: 1.284
RMSE 값: 1.296

```
1 fig = plt.figure(figsize = (25, 10))
2
3 plt.plot(my_model.predict(X_test), label='예측')
4 plt.plot(y_test.values, label='실제')
5
6 plt.legend(loc='best')
7 plt.show()
8
```



흰다리새우 가격 예측 성능평가

```

1 ## 1주 ~ 12주 전 가격을 각각 하나의 칼럼을 생성하여 저장
2 for i in range(1,13):
3     data_white_shrimp_price['{}주전 금액'.format(i)]=data_white_shrimp_price['P_PRICE'].shift(i)
4
5 rmse_list=[]

```

```

6
7 ## 2020년 이전 데이터는 Train 데이터로
8 ## 2020년 데이터는 Test 데이터로 저장
9 data_white_shrimp_20_train = data_white_shrimp_price.head(-50)
10 data_white_shrimp_20_test = data_white_shrimp_price.tail(50)
11
12 ## 2020년 가격 예측 결과를 저장할 P_PRICE_pred
13 data_white_shrimp_20_train['P_PRICE_pred'] = data_white_shrimp_20_train['P_PRICE']
14 data_white_shrimp_20_test['P_PRICE_pred'] = np.nan # 2020년 데이터의 P_PRICE_pred 칼럼은 모두 NaN
15
16 # 2020년 데이터를 한 행씩 예측 후 Train 데이터에 연결 후 연결된 새로운 데이터를 바탕으로 다시 예측
17 for i in tqdm(range(210,260)):
18     test = data_white_shrimp_20_test[data_white_shrimp_20_test.index == i]
19
20
21     X_train = data_white_shrimp_20_train.drop(['P_PRICE', 'REG_DATE', 'P_PRICE_pred'], axis=1, inplace=True)
22     y_train = data_white_shrimp_20_train['P_PRICE_pred']
23
24     X_test = data_white_shrimp_20_test.drop(['P_PRICE', 'REG_DATE', 'P_PRICE_pred'], axis=1, inplace=True)
25     y_test = data_white_shrimp_20_test['P_PRICE']
26
27     my_model = XGBRegressor(n_estimators=10000, learning_rate=0.01, seed=7857, max_depth=30, colsample_bytree=0.8)
28     my_model.fit(X_train, y_train, verbose=False)
29     pred = my_model.predict(X_test)
30     data_white_shrimp_20_test['P_PRICE_pred'] = pred # 예측 결과 값은 P_PRICE_pred 칼럼에 저장
31
32     data_white_shrimp_20_train = pd.concat([data_white_shrimp_20_train, data_white_shrimp_20_test])
33     rmse=get_rmse(y_test, pred)
34     rmse_list.append(rmse) # 각 행마다 RMSE 출력
35
36     for i in range(1,13):
37         data_white_shrimp_20_train['{}주전 금액'.format(i)]=data_white_shrimp_20_train['P_PRICE']
38
39
40 rmse_avg = np.array(rmse_list)
41 print('평균 RMSE 값은 :',rmse_avg.mean())

```

100%

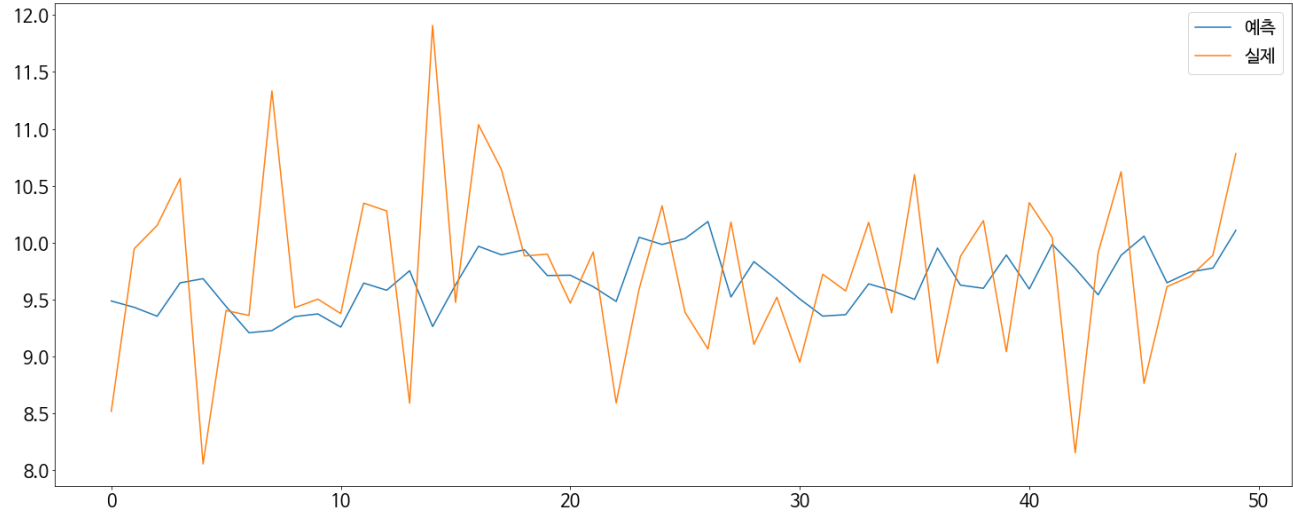
50/50 [08:30<00:00, 18.23s/it]

RMSE 값: 0.853
 RMSE 값: 0.845
 RMSE 값: 0.847
 RMSE 값: 0.849
 RMSE 값: 0.847
 RMSE 값: 0.842
 RMSE 값: 0.848
 RMSE 값: 0.847
 RMSE 값: 0.855
 RMSE 값: 0.854
 RMSE 값: 0.852
 RMSE 값: 0.853
 RMSE 값: 0.853
 RMSE 값: 0.853
 RMSE 값: 0.848
 RMSE 값: 0.846
 RMSE 값: 0.844
 RMSE 값: 0.846
 RMSE 값: 0.844
 RMSE 값: 0.845
 RMSE 값: 0.843
 RMSE 값: 0.847
 RMSE 값: 0.846
 RMSE 값: 0.847
 RMSE 값: 0.848
 RMSE 값: 0.846
 RMSE 값: 0.85
 RMSE 값: 0.846
 RMSE 값: 0.849
 RMSE 값: 0.847
 RMSE 값: 0.849
 RMSE 값: 0.851
 RMSE 값: 0.85
 RMSE 값: 0.844

```

1 fig = plt.figure(figsize = (25, 10))
2
3 plt.plot(my_model.predict(X_test), label='예측')
4 plt.plot(y_test.values, label='실제')
5
6 plt.legend(loc='best')
7 plt.show()
8

```



```
1 visualize_ftr_importances(my_model)
```

Feature Importances

원화대비 환율_MAL

▼ 7. 2021년도 예측

위의 방법을 이용해서 2020년 까지 데이터로 2021년 예측

WTi

```

1 ## 12주전 까지 데이터 생성
2
3 ## 오징어
4 for i in range(1,13):
5     data_squid_price_test['{}주전 금액'.format(i)]=data_squid_price_test['P_PRICE'].shift(i)
6
7 ## 연어
8 for i in range(1,13):
9     data_salmon_price_test['{}주전 금액'.format(i)]=data_salmon_price_test['P_PRICE'].shift(i)
10
11 ## 흰다리새우
12 for i in range(1,13):
13     data_white_shrimp_price_test['{}주전 금액'.format(i)]=data_white_shrimp_price_test['P_PRICE']

```

11주전 금액

```

1 ## 2021년은 타겟 데이터
2 data_squid_test = data_squid_price_test.iloc[260:]
3 data_salmon_test = data_salmon_price_test.iloc[260:]
4 data_white_shrimp_test = data_white_shrimp_price_test.iloc[260:]
5
6 ## 2021년 이전은 훈련 데이터
7 data_squid_train = data_squid_price_test.iloc[:260]
8 data_salmon_train = data_salmon_price_test.iloc[:260]
9 data_white_shrimp_train = data_white_shrimp_price_test.iloc[:260]
10
11
12 # 훈련데이터에서 1~12주 전의 데이터가 없는 행은 제거
13 data_squid = data_squid_train.dropna()
14 data_salmon = data_salmon_train.dropna()
15 data_white_shrimp = data_white_shrimp_train.dropna()


1 ## 2021년 오징어 가격 예측
2
3 for i in tqdm(range(260,286)):
4     test = data_squid_test[data_squid_test.index == i]
5
6     y_train = data_squid['P_PRICE']
7     X_train = data_squid.drop(['P_PRICE','REG_DATE'], axis=1, inplace=False)
8     X_test = test.drop(['P_PRICE','REG_DATE'], axis=1, inplace=False)
9
10     my_model = XGBRegressor(n_estimators=10000, learning_rate=0.001, seed=7857, max_depth=30, colsample_bytree=0.8)
11     my_model.fit(X_train, y_train, verbose=False)
12     pred = my_model.predict(X_test)
13     test['P_PRICE'] = pred
14

```



```

15 data_squid = pd.concat([data_squid, test])
16
17 for i in range(1,13):
18     data_squid['{}주전 금액'.format(i)]=data_squid['P_PRICE'].shift(i)
19

```

100%

26/26 [02:57<00:00, 7.09s/it]

```

1 ## 2021년 연어 가격 예측
2
3 for i in tqdm(range(262,288)):
4     test = data_salmon_test[data_salmon_test.index == i]
5
6     y_train = data_salmon['P_PRICE']
7     X_train = data_salmon.drop(['P_PRICE','REG_DATE'], axis=1, inplace=False)
8     X_test = test.drop(['P_PRICE','REG_DATE'], axis=1, inplace=False)
9
10
11     my_model = XGBRegressor(n_estimators=10000, learning_rate=0.001, seed=7857, max_depth=30, colsamplebytree=0.8)
12     my_model.fit(X_train, y_train, verbose=False)
13     pred = my_model.predict(X_test)
14     test['P_PRICE'] = pred
15
16     data_salmon = pd.concat([data_salmon, test])
17
18     for i in range(1,13):
19         data_salmon['{}주전 금액'.format(i)]=data_salmon['P_PRICE'].shift(i)
20

```

100%

26/26 [02:45<00:00, 6.64s/it]

```

1 ## 2021년 흰다리새우 가격 예측
2
3 for i in tqdm(range(262,288)):
4     test = data_white_shrimp_test[data_white_shrimp_test.index == i]
5
6     y_train = data_white_shrimp['P_PRICE']
7     X_train = data_white_shrimp.drop(['P_PRICE','REG_DATE'], axis=1, inplace=False)
8     X_test = test.drop(['P_PRICE','REG_DATE'], axis=1, inplace=False)
9
10
11     my_model = XGBRegressor(n_estimators=10000, learning_rate=0.001, seed=7857, max_depth=30, colsamplebytree=0.8)
12     my_model.fit(X_train, y_train, verbose=False)
13     pred = my_model.predict(X_test)
14     test['P_PRICE'] = pred
15
16     data_white_shrimp = pd.concat([data_white_shrimp, test])
17
18     for i in range(1,13):
19         data_white_shrimp['{}주전 금액'.format(i)]=data_white_shrimp['P_PRICE'].shift(i)
20

```

100%

26/26 [03:03<00:00, 7.53s/it]

▼ 8. 최종결과물

```
1 sub_data_squid=data_squid[['REG_DATE','P_PRICE']].tail(26)
2 sub_data_squid.reset_index(drop=True,inplace=True)
3 sub_data_squid.to_excel('/content/drive/MyDrive/TAVE_빅콘테스트/제출데이터/2021_오징어_단가예측
```

```
1 sub_data_salmon=data_salmon[['REG_DATE','P_PRICE']].tail(26)
2 sub_data_salmon.reset_index(drop=True,inplace=True)
3 sub_data_salmon.to_excel('/content/drive/MyDrive/TAVE_빅콘테스트/제출데이터/2021_연어_단가예측.;
```

```
1 sub_data_white_shrimp=data_white_shrimp[['REG_DATE','P_PRICE']].tail(26)
2 sub_data_white_shrimp.reset_index(drop=True,inplace=True)
3 sub_data_white_shrimp.to_excel('/content/drive/MyDrive/TAVE_빅콘테스트/제출데이터/2021_흰다리새;
```

1

✓ 0초 오후 10:58에 완료됨

● ✕