

```
In [184]: import numpy as np
import pandas as pd

import matplotlib as mpl
import seaborn as sns
plt.rcParams['font.family'] = 'sans-serif'
plt.rcParams['axes.uncheckedaxis'] = True
```

```
In [184]: sw = pd.read_csv('C:\Users\7hds\Downloads\sw.csv', encoding='cp949')
```

서울교통공사에서 제공한 시간별 지하철 혼잡도 데이터. 일주일이었어. 일주일이었어.

```
In [187]: sw.sum()

Out[187]:
time      0
종이값    0
역번호    0
역명      0
구분      0
613000    0
613000    0
713000    0
713000    0
813000    0
813000    0
913000    0
913000    0
1013000   0
1013000   0
1113000   0
1113000   0
1213000   0
1213000   0
1313000   0
1313000   0
1413000   0
1413000   0
1513000   0
1513000   0
1613000   0
1613000   0
1713000   0
1713000   0
1813000   0
1813000   0
1913000   0
1913000   0
2013000   0
2013000   0
2113000   0
2113000   0
2213000   0
2213000   0
2313000   0
2313000   0
dtype: int64

사용할 연속형 변수들만 간추린 데이터셋 생성.
```

```
In [188]: sw2 = sw.drop(['time', '종이값', '역번호', '역명', '구분'], axis = 1)
```

```
In [188]: sw2.info()

Out[188]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1784 entries, 0 to 1783
Data columns (total 7 columns):
#   column  Non-Null Count  Dtype
---  --
0   1813000  1784 non-null   float64
1   613000   1784 non-null   float64
2   613000   1784 non-null   float64
3   713000   1784 non-null   float64
4   713000   1784 non-null   float64
5   813000   1784 non-null   float64
6   813000   1784 non-null   float64
7   913000   1784 non-null   float64
8   913000   1784 non-null   float64
9   1013000  1784 non-null  float64
10  1013000  1784 non-null  float64
11  1113000  1784 non-null  float64
12  1113000  1784 non-null  float64
13  1213000  1784 non-null  float64
14  1213000  1784 non-null  float64
15  1313000  1784 non-null  float64
16  1313000  1784 non-null  float64
17  1413000  1784 non-null  float64
18  1413000  1784 non-null  float64
19  1513000  1784 non-null  float64
20  1513000  1784 non-null  float64
21  1613000  1784 non-null  float64
22  1613000  1784 non-null  float64
23  1713000  1784 non-null  float64
24  1713000  1784 non-null  float64
25  1813000  1784 non-null  float64
26  1813000  1784 non-null  float64
27  1913000  1784 non-null  float64
28  1913000  1784 non-null  float64
29  2013000  1784 non-null  float64
30  2013000  1784 non-null  float64
31  2113000  1784 non-null  float64
32  2113000  1784 non-null  float64
33  2213000  1784 non-null  float64
34  2213000  1784 non-null  float64
35  2313000  1784 non-null  float64
36  2313000  1784 non-null  float64
dtypes: float64(37)
memory usage: 492.7 KB

QR방식을 사용, 이상치를 제거함.
```

```
In [188]: quartile_3 = sw2.quantile(0.75)
quartile_3 = sw2.quantile(0.75)
condition = (sw2 < quartile_3 - 1.5 * IQR) | (sw2 > (quartile_3 + 1.5 * IQR))
condition = condition.any(axis=1)
sw2 = sw2[condition]
```

각 시간별 혼잡도의 평균을 구해 본 결과 혼잡도가 가장 낮은 23시 30분과 혼잡도가 가장 높은 18시 30분을 중심으로 분석 진행.

```
In [189]: sw2.mean()

Out[189]:
1813000    21.799549
613000     22.864293
613000     24.599995
713000     22.864293
713000     21.853085
813000     27.891895
813000     27.899995
913000     42.899999
913000     42.899999
1013000    29.842980
1013000    29.899999
1113000    28.749774
1113000    28.829444
1213000    29.710623
1213000    29.842980
1313000    31.842980
1313000    31.499981
1413000    29.891121
1413000    30.779988
1513000    32.371493
1513000    32.292974
1613000    30.949929
1613000    30.779988
1713000    43.495439
1713000    43.499999
1813000    51.945159
1813000    51.942944
1913000    38.184751
1913000    29.849944
2013000    26.313389
2013000    25.779988
2113000    27.849999
2113000    25.742925
2213000    30.897959
2213000    29.897944
2313000    24.907249
2313000    17.889995
dtype: float64

Kmeans 군집화를 하기 앞서 변수를 정규화시킴.
```

```
In [189]: from sklearn.preprocessing import MinMaxScaler

In [189]: #혼잡도 평균이 가장 낮은 시간인 23시 30분과 혼잡도 평균이 가장 높은 18시 30분
data = sw2[['1813000', '2313000']]

# scaler 생성
scaler = MinMaxScaler()
data_scale = scaler.fit_transform(data)
```

```
In [191]: from sklearn.cluster import KMeans

k = 2

# 그림을 그릴 때, random_state 설정
model = KMeans(n_clusters = k, random_state = 10)
model.fit(data_scale)

# 군집을 할 데이터에 대해
model.fit(data_scale)

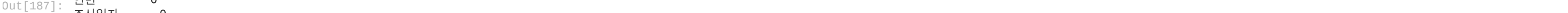
# 클러스터링 결과 각 데이터가 몇 번째 군집에 속하는지 저장
sw2['cluster'] = model.fit_predict(data_scale)

import matplotlib.pyplot as plt
plt.figure(figsize = (8, 8))

for i in range(k):
    plt.scatter(sw2.loc[sw2['cluster'] == i, '1813000'], sw2.loc[sw2['cluster'] == i, '2313000'],
                label = 'cluster ' + str(i))

plt.legend()
plt.title('k = %d results'%k, size = 15)
plt.xlabel('1813000', size = 12)
plt.ylabel('2313000', size = 12)
plt.show()

from yellowbrick.cluster import SilhouetteVisualizer
visualizer = SilhouetteVisualizer(model, colors='yellowbrick')
visualizer.fit(data_scale)
visualizer.show()
```



```
Out[192]: <AxesSubplot: title='center': 'Silhouette Plot of KMeans Clustering for 442 Samples in 2 Centers', xlabel='silhouette coefficient values', ylabel='cluster label'>
```

```
In [192]: k = 3

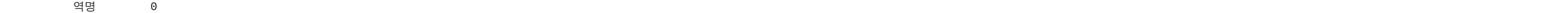
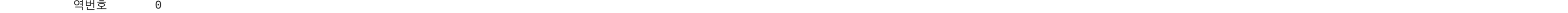
model = KMeans(n_clusters = k, random_state = 10)
model.fit(data_scale)
sw2['cluster'] = model.fit_predict(data_scale)

import matplotlib.pyplot as plt
plt.figure(figsize = (8, 8))

for i in range(k):
    plt.scatter(sw2.loc[sw2['cluster'] == i, '1813000'], sw2.loc[sw2['cluster'] == i, '2313000'],
                label = 'cluster ' + str(i))

plt.legend()
plt.title('k = %d results'%k, size = 15)
plt.xlabel('1813000', size = 12)
plt.ylabel('2313000', size = 12)
plt.show()

from yellowbrick.cluster import SilhouetteVisualizer
visualizer_2 = SilhouetteVisualizer(model, colors='yellowbrick')
visualizer_2.fit(data_scale)
visualizer_2.show()
```



```
Out[193]: <AxesSubplot: title='center': 'Silhouette Plot of KMeans Clustering for 442 Samples in 3 Centers', xlabel='silhouette coefficient values', ylabel='cluster label'>
```

```
In [193]: k = 4

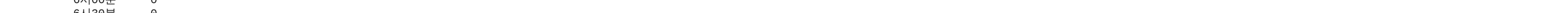
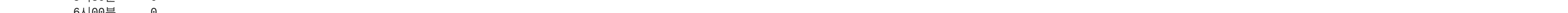
model = KMeans(n_clusters = k, random_state = 10)
model.fit(data_scale)
sw2['cluster'] = model.fit_predict(data_scale)

import matplotlib.pyplot as plt
plt.figure(figsize = (8, 8))

for i in range(k):
    plt.scatter(sw2.loc[sw2['cluster'] == i, '1813000'], sw2.loc[sw2['cluster'] == i, '2313000'],
                label = 'cluster ' + str(i))

plt.legend()
plt.title('k = %d results'%k, size = 15)
plt.xlabel('1813000', size = 12)
plt.ylabel('2313000', size = 12)
plt.show()

from yellowbrick.cluster import SilhouetteVisualizer
visualizer_2 = SilhouetteVisualizer(model, colors='yellowbrick')
visualizer_2.fit(data_scale)
visualizer_2.show()
```



```
Out[194]: <AxesSubplot: title='center': 'Silhouette Plot of KMeans Clustering for 442 Samples in 4 Centers', xlabel='silhouette coefficient values', ylabel='cluster label'>
```

```
In [194]: k = 5

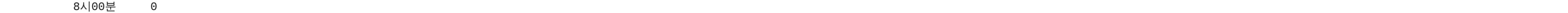
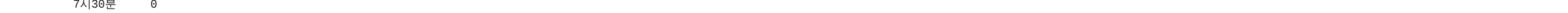
model = KMeans(n_clusters = k, random_state = 10)
model.fit(data_scale)
sw2['cluster'] = model.fit_predict(data_scale)

import matplotlib.pyplot as plt
plt.figure(figsize = (8, 8))

for i in range(k):
    plt.scatter(sw2.loc[sw2['cluster'] == i, '1813000'], sw2.loc[sw2['cluster'] == i, '2313000'],
                label = 'cluster ' + str(i))

plt.legend()
plt.title('k = %d results'%k, size = 15)
plt.xlabel('1813000', size = 12)
plt.ylabel('2313000', size = 12)
plt.show()

from yellowbrick.cluster import SilhouetteVisualizer
visualizer_2 = SilhouetteVisualizer(model, colors='yellowbrick')
visualizer_2.fit(data_scale)
visualizer_2.show()
```



```
Out[195]: <AxesSubplot: title='center': 'Silhouette Plot of KMeans Clustering for 442 Samples in 5 Centers', xlabel='silhouette coefficient values', ylabel='cluster label'>
```

```
In [195]: k = 6

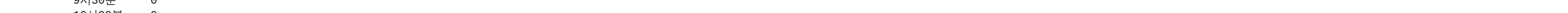
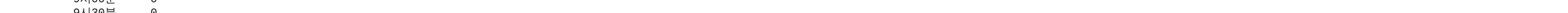
model = KMeans(n_clusters = k, random_state = 10)
model.fit(data_scale)
sw2['cluster'] = model.fit_predict(data_scale)

import matplotlib.pyplot as plt
plt.figure(figsize = (8, 8))

for i in range(k):
    plt.scatter(sw2.loc[sw2['cluster'] == i, '1813000'], sw2.loc[sw2['cluster'] == i, '2313000'],
                label = 'cluster ' + str(i))

plt.legend()
plt.title('k = %d results'%k, size = 15)
plt.xlabel('1813000', size = 12)
plt.ylabel('2313000', size = 12)
plt.show()

from yellowbrick.cluster import SilhouetteVisualizer
visualizer_2 = SilhouetteVisualizer(model, colors='yellowbrick')
visualizer_2.fit(data_scale)
visualizer_2.show()
```



```
Out[196]: <AxesSubplot: title='center': 'Silhouette Plot of KMeans Clustering for 442 Samples in 6 Centers', xlabel='silhouette coefficient values', ylabel='cluster label'>
```

시각화 결과를 얻으면서도 보았을 때, 최적의 K는 3~5사이로 보임.

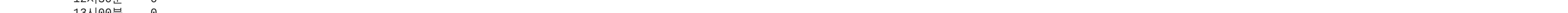
또한 Cluster 단의 개리와 값을 나타내는 값소 정도를 K를 뜻하여 시각화 시각화 결과 최적의 K값은 3이라고 보임.

```
In [197]: from yellowbrick.cluster import KElbowVisualizer

model = KMeans()
visualizer = KElbowVisualizer(model, k=(1,10))
visualizer.fit(data_scale)
visualizer.show()
```

D:\V\Dev\test-package>sklearn\cluster_kmeans.py:1038: UserWarning: KMeans is known to have a memory leak on Windows with MSVC, when there are less chunks than available threads. You can avoid it by setting the environment variable env.NUMPY_THREADS=2.

```
Out[197]: KElbowVisualizer(ax=AxesSubplot(), estimator=KMeans(n_clusters=6), k=(1, 10))
```



In []: