

## 20장. 군집화 (Clustering)

In [1]:

```
1 !pip install seaborn
```

Requirement already satisfied: seaborn in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (0.11.1)  
 Requirement already satisfied: pandas>=0.23 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from seaborn) (1.2.3)  
 Requirement already satisfied: numpy>=1.15 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from seaborn) (1.20.1)  
 Requirement already satisfied: matplotlib>=2.2 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from seaborn) (3.3.4)  
 Requirement already satisfied: scipy>=1.0 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from seaborn) (1.6.3)  
 Requirement already satisfied: cycler>=0.10 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from matplotlib>=2.2->seaborn) (0.10.0)  
 Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.3 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from matplotlib>=2.2->seaborn) (2.4.7)  
 Requirement already satisfied: pillow>=6.2.0 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from matplotlib>=2.2->seaborn) (8.1.2)  
 Requirement already satisfied: python-dateutil>=2.1 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from matplotlib>=2.2->seaborn) (2.8.1)  
 Requirement already satisfied: kiwisolver>=1.0.1 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from matplotlib>=2.2->seaborn) (1.3.1)  
 Requirement already satisfied: six in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from cycler>=0.10->matplotlib>=2.2->seaborn) (1.15.0)  
 Requirement already satisfied: pytz>=2017.3 in c:\Users\WghkrG\Anaconda3\envs\data\_mining\lib\site-packages (from pandas>=0.23->seaborn) (2021.1)

### 1. 데이터셋 (IRIS)

<https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data> (<https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data>)

In [2]:

```
1 import requests
2 import os
3
4 data = requests.get("https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data")
5 path = os.path.join('data', 'iris.data')
6 with open(path, "w") as f:
7     f.write(data.text)
```

#### 1.1 데이터셋 읽기

In [3]:

```

1 import pandas as pd
2 column_names = ['sepal length', 'sepal width', 'petal length', 'petal width', 'species']
3 dataset = pd.read_csv(path, names=column_names)
4 dataset.head()

```

Out[3]:

	sepal length	sepal width	petal length	petal width	species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

In [4]:

```
1 dataset.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0   sepal length    150 non-null   float64
1   sepal width     150 non-null   float64
2   petal length    150 non-null   float64
3   petal width     150 non-null   float64
4   species         150 non-null   object
dtypes: float64(4), object(1)
memory usage: 6.0+ KB

```

## 2. 데이터 탐색

### 2.1 요약 통계량

In [5]:

```
1 dataset.describe()
```

Out[5]:

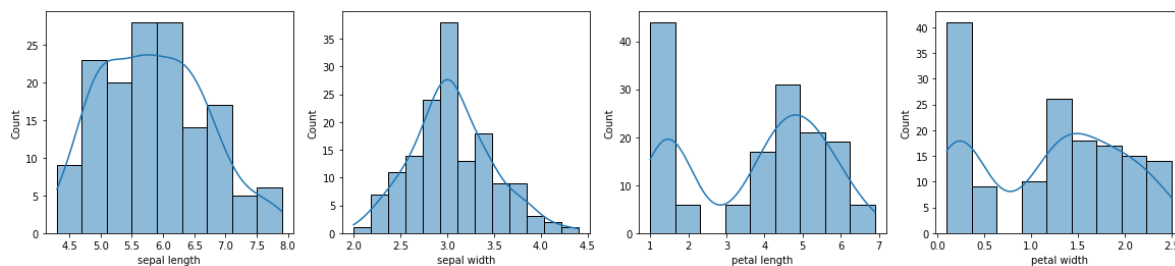
	sepal length	sepal width	petal length	petal width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

## 2.2 단일 변수 분석

### 2.2.1 히스토그램

In [6]:

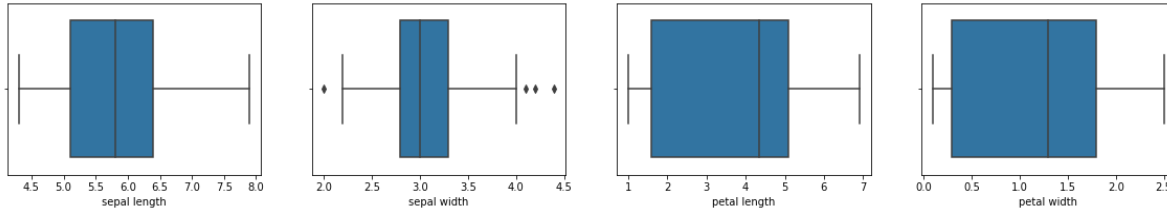
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 plt.figure(figsize=[20,4])
5 for i, column in enumerate(dataset.describe().columns):
6     plt.subplot(1,4,i+1)
7     sns.histplot(data=dataset, x=column, kde=True)
8 plt.show()
```



### 2.2.2 박스 플롯

In [7]:

```
1 plt.figure(figsize=[20,3])
2 for i in enumerate(dataset.describe().columns[:4]):
3     plt.subplot(1,4,i[0]+1)
4     sns.boxplot(x=dataset[i[1]])
5
6 plt.show()
```

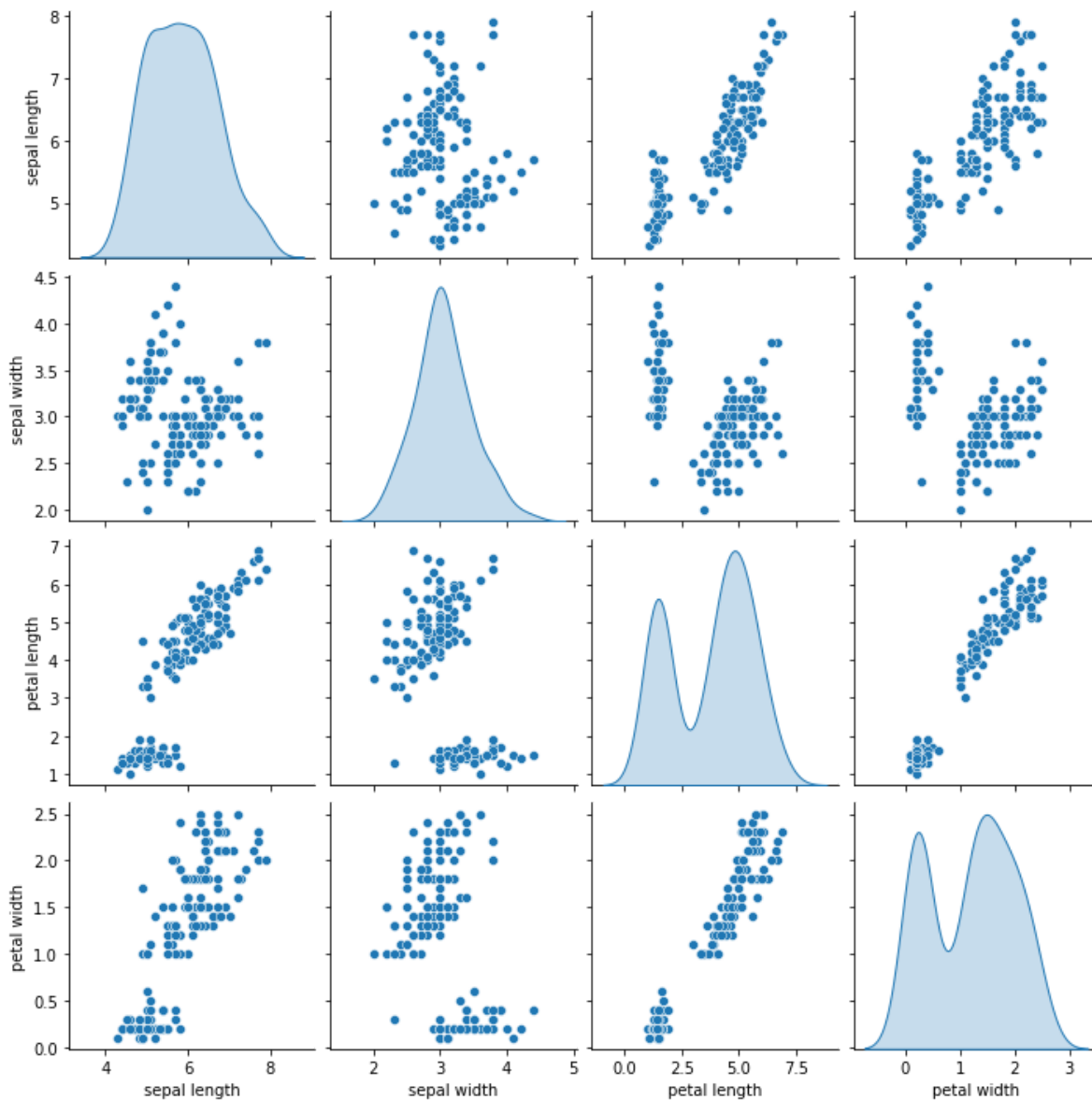


## 2.3 두 변수 관계 분석

### 2.3.1 산포도 행렬

In [8]:

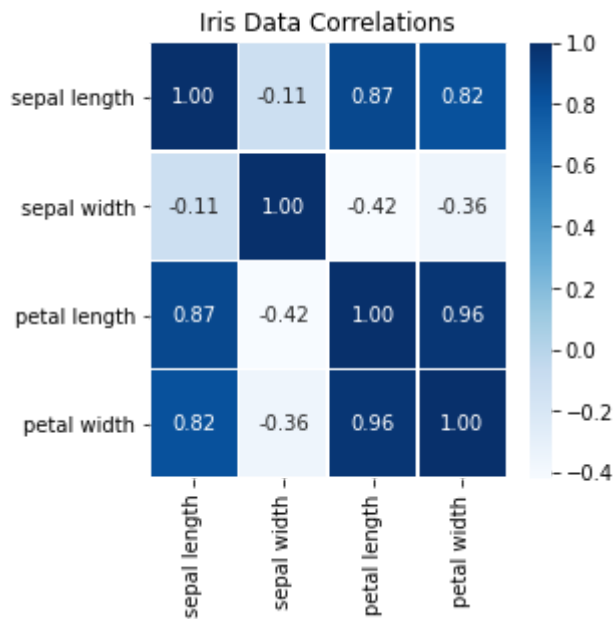
```
1 sns.pairplot(dataset, diag_kind="kde")  
2 plt.show()
```



### 2.3.2 히트맵

In [9]:

```
1 fig, ax = plt.subplots(figsize=(4, 4))
2 sns.heatmap(dataset.corr(), linewidths=.5, annot=True, fmt=".2f", cmap='Blues')
3 plt.title('Iris Data Correlations')
4 plt.show()
```



### 3. 데이터 전처리

#### 3.1 데이터 추출

In [10]:

```

1 clusterdata = dataset.iloc[:, :-1]
2 inputs = clusterdata.iloc[:, :].values.tolist()
3 columns = clusterdata.keys().tolist()
4 column2index = {column : i for i, column in enumerate(columns)}
5 print('columns = ', columns)
6 print('column2index = ', column2index)

```

```

columns = ['sepal length', 'sepal width', 'petal length', 'petal width']
column2index = {'sepal length': 0, 'sepal width': 1, 'petal length': 2, 'petal width': 3}

```

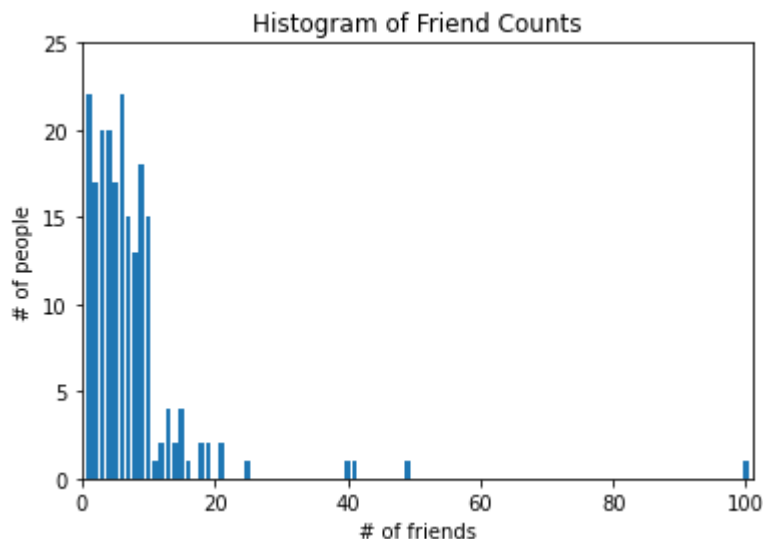
## 3.2 데이터 표준화

In [11]:

```

1 from scratch.working_with_data import scale, rescale, Vector
2 from typing import List
3
4 inputs_normed = rescale(inputs)

```



## 4. K-평균 군집화

### Q1. 손실 곡선을 보고 K 선택하기

손실을 최소화 하는 클러스터 수  $K$ 를 찾아보시오. 단,  $K$ 는 20까지 확인해 보라.

In [12]:

```

1 from scratch.linear_algebra import Vector
2
3 def num_differences(v1: Vector, v2: Vector) -> int:
4     assert len(v1) == len(v2)
5     return len([x1 for x1, x2 in zip(v1, v2) if x1 != x2])

```

In [13]:

```
1 from typing import List
2 from scratch.linear_algebra import vector_mean
3
4 def cluster_means(k: int,
5                   inputs: List[Vector],
6                   assignments: List[int]) -> List[Vector]:
7     clusters = [[] for i in range(k)]
8     for input, assignment in zip(inputs, assignments):
9         clusters[assignment].append(input)
10
11     return [vector_mean(cluster) if cluster else random.choice(inputs) for cluster in clusters]
```

In [14]:

```
1 import itertools
2 import random
3 import tqdm
4 from scratch.linear_algebra import squared_distance
5
6 class KMeans:
7     def __init__(self, k: int) -> None:
8         self.k = k
9         self.means = None
10
11     def classify(self, input: Vector) -> int:
12         return min(range(self.k), key=lambda i: squared_distance(input, self.means[i]))
13
14     def train(self, inputs: List[Vector]) -> None:
15         assignments = [random.randrange(self.k) for _ in inputs]
16         with tqdm.tqdm(itertools.count()) as t:
17             for _ in t:
18                 self.means = cluster_means(self.k, inputs, assignments)
19                 new_assignments = [self.classify(input) for input in inputs]
20                 num_changed = num_differences(assignments, new_assignments)
21                 if num_changed == 0:
22                     return
23
24                 assignments = new_assignments
25                 t.set_description(f"changed: {num_changed} / {len(inputs)}")
```



In [15]:

```

1  # your code
2  import random
3  from matplotlib import pyplot as plt
4
5  def squared_clustering_errors(inputs: List[Vector], k: int) -> float:
6      clusterer = KMeans(k)
7      clusterer.train(inputs)
8      means = clusterer.means
9      assignments = [clusterer.classify(input) for input in inputs]
10
11     return sum(squared_distance(input, means[cluster]) for input, cluster in zip(inputs, assignments))
12
13
14  ks = range(1, 21)
15  errors = [squared_clustering_errors(inputs, k) for k in ks]
16
17  optimal_k = errors.index(min(errors)) + 1
18
19  fig, ax = plt.subplots(figsize=(10, 5))
20  plt.plot(ks, errors)
21  plt.xticks(ks)
22  plt.xlabel(f"k (Optimal k = {optimal_k})")
23  plt.ylabel("total squared error")
24  plt.title("Total Error vs. # of Clusters")
25  plt.show()

```

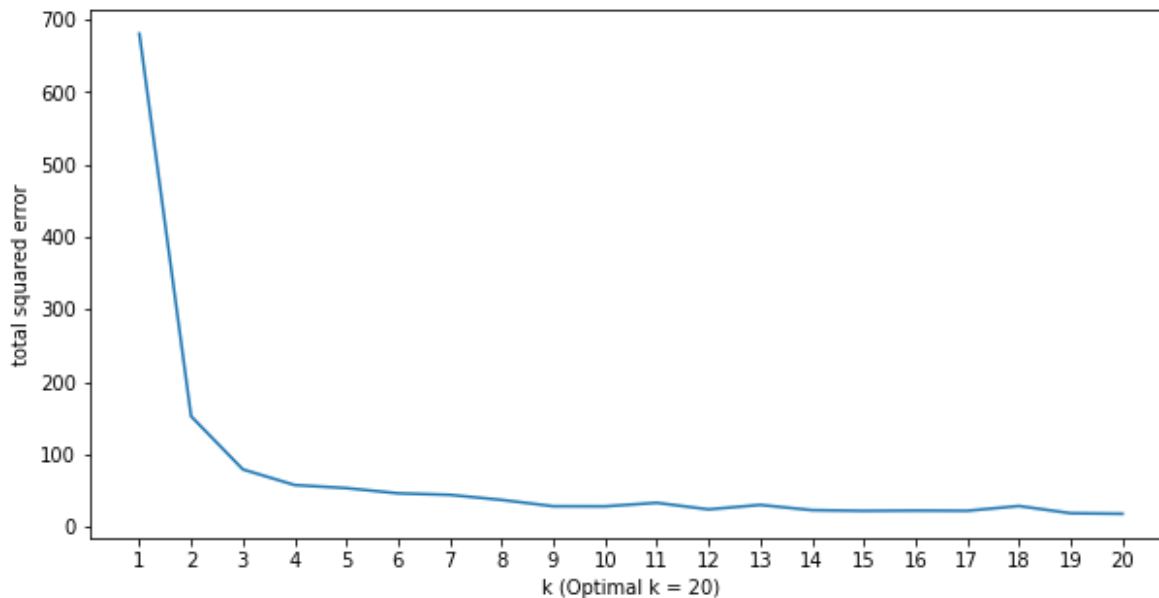
Out [00:00, ?it/s]

```

changed: 1 / 150: : 3it [00:00, 749.74it/s]
changed: 1 / 150: : 10it [00:00, 769.05it/s]
changed: 1 / 150: : 8it [00:00, 615.26it/s]
changed: 1 / 150: : 6it [00:00, 545.35it/s]
changed: 2 / 150: : 7it [00:00, 499.91it/s]
changed: 3 / 150: : 11it [00:00, 439.91it/s]
changed: 4 / 150: : 9it [00:00, 374.92it/s]
changed: 1 / 150: : 6it [00:00, 315.70it/s]
changed: 1 / 150: : 9it [00:00, 310.27it/s]
changed: 2 / 150: : 8it [00:00, 296.23it/s]
changed: 3 / 150: : 12it [00:00, 272.66it/s]
changed: 2 / 150: : 5it [00:00, 217.34it/s]
changed: 4 / 150: : 6it [00:00, 230.72it/s]
changed: 1 / 150: : 9it [00:00, 224.95it/s]
changed: 3 / 150: : 6it [00:00, 199.96it/s]
changed: 1 / 150: : 6it [00:00, 181.77it/s]
changed: 1 / 150: : 6it [00:00, 181.78it/s]
changed: 1 / 150: : 7it [00:00, 174.96it/s]
changed: 1 / 150: : 14it [00:00, 177.18it/s]

```

Total Error vs. # of Clusters



## Q2. 군집화 및 결과 확인 (Q)

$K=3$ 으로 군집화를 해서 다음과 같이 군집화 결과를 확인해 보라.

In [16]:

```
1 from typing import NamedTuple, Union
2
3 class Leaf(NamedTuple):
4     value: Vector
```

In [17]:

```
1 class Merged(NamedTuple):
2     children: tuple
3     order: int
```

In [18]:

```
1 Cluster = Union[Leaf, Merged]
```

In [19]:

```
1 def get_values(cluster: Cluster) -> List[Vector]:
2     if isinstance(cluster, Leaf):
3         return [cluster.value]
4     else:
5         return [value
6                 for child in cluster.children
7                 for value in get_values(child)]
```

In [20]:

```
1 from typing import Callable
2 from scratch.linear_algebra import distance
3
4 def cluster_distance(cluster1: Cluster,
5                       cluster2: Cluster,
6                       distance_agg: Callable = min) -> float:
7     return distance_agg([distance(v1, v2)
8                           for v1 in get_values(cluster1)
9                           for v2 in get_values(cluster2)])
```

In [21]:

```
1 def get_merge_order(cluster: Cluster) -> float:
2     if isinstance(cluster, Leaf):
3         return float('inf')
4     else:
5         return cluster.order
```

In [22]:

```
1 from typing import Tuple
2
3 def get_children(cluster: Cluster):
4     if isinstance(cluster, Leaf):
5         raise TypeError("Leaf has no children")
6     else:
7         return cluster.children
```

In [23]:

```
1 def bottom_up_cluster(inputs: List[Vector],
2                       distance_agg: Callable = min) -> Cluster:
3     clusters: List[Cluster] = [Leaf(input) for input in inputs]
```

In [24]:

```
1 def pair_distance(pair: Tuple[Cluster, Cluster]) -> float:
2     return cluster_distance(pair[0], pair[1], distance_agg)
```

In [25]:

```
1 def generate_clusters(base_cluster: Cluster,
2                       num_clusters: int) -> List[Cluster]:
3     clusters = [base_cluster]
4
5     while len(clusters) < num_clusters:
6         next_cluster = min(clusters, key=get_merge_order)
7         clusters = [c for c in clusters if c != next_cluster]
8         clusters.extend(get_children(next_cluster))
9     return clusters
```

In [26]:

```

1 #your code
2 #base_cluster = bottom_up_cluster(inputs)
3 #three_clusters = [get_values(cluster)
4 #                   for cluster in generate_clusters(base_cluster, 3)]
5
6 random.seed(0)
7 clusterer = KMeans(k=3)
8 clusterer.train(inputs)
9 means = sorted(clusterer.means)
10 assignments = [clusterer.classify(input) for input in inputs]
11
12 assert len(means) == 3

```

chagned: 1 / 150: : 11it [00:00, 785.53it/s]

**dataset에 k\_means 컬럼 추가**

In [27]:

```

1 dataset["k_means"] = assignments
2 dataset.head()

```

Out[27]:

	sepal length	sepal width	petal length	petal width	species	k_means
0	5.1	3.5	1.4	0.2	Iris-setosa	1
1	4.9	3.0	1.4	0.2	Iris-setosa	1
2	4.7	3.2	1.3	0.2	Iris-setosa	1
3	4.6	3.1	1.5	0.2	Iris-setosa	1
4	5.0	3.6	1.4	0.2	Iris-setosa	1

**species와 k\_means 결과 비교**

In [28]:

```
1 dataset[dataset['k_means']!=0].head()
```

Out[28]:

	sepal length	sepal width	petal length	petal width	species	k_means
50	7.0	3.2	4.7	1.4	Iris-versicolor	0
52	6.9	3.1	4.9	1.5	Iris-versicolor	0
77	6.7	3.0	5.0	1.7	Iris-versicolor	0
100	6.3	3.3	6.0	2.5	Iris-virginica	0
102	7.1	3.0	5.9	2.1	Iris-virginica	0

In [29]:

```
1 dataset.groupby(["k_means", "species"])["k_means"].count()
```

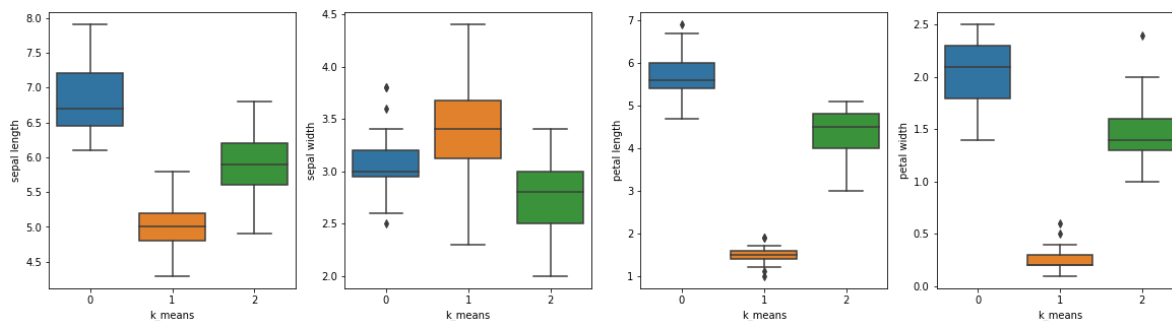
Out[29]:

```
k_means  species
0        Iris-versicolor    3
         Iris-virginica    36
1        Iris-setosa       50
2        Iris-versicolor   47
         Iris-virginica    14
Name: k_means, dtype: int64
```

## 군집화 결과 시각화

In [30]:

```
1 plt.subplots(figsize=(20, 5))
2 plt.subplot(1,4,1)
3 sns.boxplot(x='k_means', y='sepal length', data=dataset)
4 plt.subplot(1,4,2)
5 sns.boxplot(x='k_means', y='sepal width', data=dataset)
6 plt.subplot(1,4,3)
7 sns.boxplot(x='k_means', y='petal length', data=dataset)
8 plt.subplot(1,4,4)
9 sns.boxplot(x='k_means', y='petal width', data=dataset)
10 plt.show()
```



## Q3. 군집화 및 결과 확인

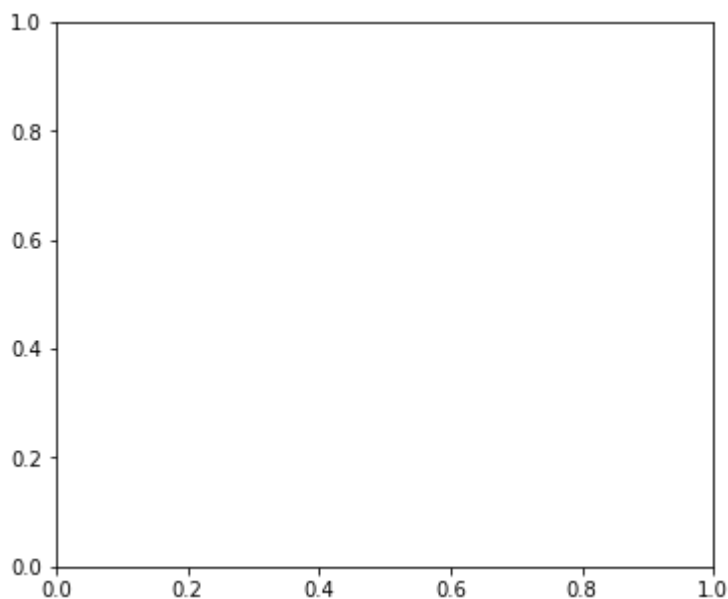
각 군집이 구분되도록 두 변수의 산포도를 그리는 함수 `plot_cluster` 구현하시오.

In [31]:

```
1 plt.subplots(figsize=(20, 5))
2 plt.subplot(1,3,1)
3 plot_cluster(clusters, column2index["sepal length"], column2index["sepal width"])
4 plt.subplot(1,3,2)
5 plot_cluster(clusters, column2index["petal length"], column2index["petal width"])
6 plt.subplot(1,3,3)
7 plot_cluster(clusters, column2index["sepal length"], column2index["petal length"])
8 plt.show()
```

```
-----
-
NameError                                Traceback (most recent call last)
<ipython-input-31-6fa796c3a7a0> in <module>
      1 plt.subplots(figsize=(20, 5))
      2 plt.subplot(1,3,1)
----> 3 plot_cluster(clusters, column2index["sepal length"], column2index["sepal
width"])
      4 plt.subplot(1,3,2)
      5 plot_cluster(clusters, column2index["petal length"], column2index["petal
width"])
```

NameError: name 'plot\_cluster' is not defined



## 5. 계층 군집화 (Hierarchical Clustering)

### Q4. 군집화 및 결과 확인

$K=3$ 으로 최장 거리(max) 기준으로 군집화를 해서 다음과 같이 군집화 결과를 확인해 보라.

In [32]:

```
1 # your code
```

**dataset에 h\_clustering 컬럼 추가**

In [33]:

```
1 dataset["h_clustering"] = h_assignments
2 dataset.head()
```

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-33-4f021a770fb5> in <module>
----> 1 dataset["h_clustering"] = h_assignments
      2 dataset.head()

NameError: name 'h_assignments' is not defined
```

**h\_clustering과 species 비교**

In [34]:

```
1 dataset.groupby(["h_clustering", "species"])['h_clustering'].count()
```

```
-----
-
KeyError                                Traceback (most recent call last)
<ipython-input-34-534101fc94bb> in <module>
----> 1 dataset.groupby(["h_clustering", "species"])['h_clustering'].count()

~Wanaconda3WenvsWdata_mininingWlibWsite-packagesWpandasWcoreWframe.py in grou
pby(self, by, axis, level, as_index, sort, group_keys, squeeze, observed,
dropna)
    6715         axis = self._get_axis_number(axis)
    6716
-> 6717         return DataFrameGroupBy(
    6718             obj=self,
    6719             keys=by,

~Wanaconda3WenvsWdata_mininingWlibWsite-packagesWpandasWcoreWgroupbyWgroupb
y.py in __init__(self, obj, keys, axis, level, grouper, exclusions, selecti
on, as_index, sort, group_keys, squeeze, observed, mutated, dropna)
    558         from pandas.core.groupby.grouper import get_grouper
    559
--> 560         grouper, exclusions, obj = get_grouper(
    561             obj,
    562             keys,

~Wanaconda3WenvsWdata_mininingWlibWsite-packagesWpandasWcoreWgroupbyWgroupe
r.py in get_grouper(obj, key, axis, level, sort, observed, mutated, validat
e, dropna)
    809         in_axis, name, level, gpr = False, None, gpr, None
    810     else:
--> 811         raise KeyError(gpr)
    812     elif isinstance(gpr, Grouper) and gpr.key is not None:
    813         # Add key to exclusions

KeyError: 'h_clustering'
```

## 군집화 결과 시각화



In [35]:

```

1 plt.subplots(figsize=(20, 5))
2 plt.subplot(1,4,1)
3 sns.boxplot(x = 'h_clustering', y = 'sepal length', data= dataset)
4 plt.subplot(1,4,2)
5 sns.boxplot(x = 'h_clustering', y = 'sepal width', data= dataset)
6 plt.subplot(1,4,3)
7 sns.boxplot(x = 'h_clustering', y = 'petal length', data= dataset)
8 plt.subplot(1,4,4)
9 sns.boxplot(x = 'h_clustering', y = 'petal width', data= dataset)
10 plt.show()

```

```

-----
ValueError                                Traceback (most recent call last)
<ipython-input-35-49390ab93351> in <module>
      1 plt.subplots(figsize=(20, 5))
      2 plt.subplot(1,4,1)
----> 3 sns.boxplot(x = 'h_clustering', y = 'sepal length', data= dataset)
      4 plt.subplot(1,4,2)
      5 sns.boxplot(x = 'h_clustering', y = 'sepal width', data= dataset)

~Wanaconda3WenvsWdata_miningWlibWsite-packagesWseabornWdecorators.py in i
nner_f(*args, **kwargs)
     44     )
     45     kwargs.update({k: arg for k, arg in zip(sig.parameters, args)})
--> 46     return f(**kwargs)
     47     return inner_f
     48

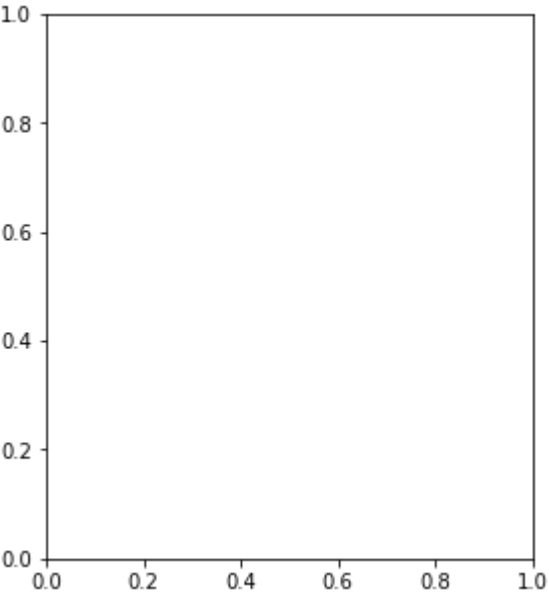
~Wanaconda3WenvsWdata_miningWlibWsite-packagesWseabornWcategorical.py in b
oxplot(x, y, hue, data, order, hue_order, orient, color, palette, saturati
on, width, dodge, fliersize, linewidth, whis, ax, **kwargs)
    2238 ):
    2239
-> 2240     plotter = _BoxPlotter(x, y, hue, data, order, hue_order,
    2241                             orient, color, palette, saturation,
    2242                             width, dodge, fliersize, linewidth)

~Wanaconda3WenvsWdata_miningWlibWsite-packagesWseabornWcategorical.py in _
_init__(self, x, y, hue, data, order, hue_order, orient, color, palette, s
aturation, width, dodge, fliersize, linewidth)
    404         width, dodge, fliersize, linewidth):
    405
--> 406         self.establish_variables(x, y, hue, data, orient, order, hue_orde
r)
    407         self.establish_colors(color, palette, saturation)
    408

~Wanaconda3WenvsWdata_miningWlibWsite-packagesWseabornWcategorical.py in e
stablish_variables(self, x, y, hue, data, orient, order, hue_order, units)
    151         if isinstance(var, str):
    152             err = "Could not interpret input '{}'.format(var)
--> 153             raise ValueError(err)
    154
    155         # Figure out the plotting orientation

ValueError: Could not interpret input 'h_clustering'

```



In [36]:

```
1 plt.subplots(figsize=(20, 5))
2 plt.subplot(1,3,1)
3 plot_cluster(h_clusters, column2index["sepal length"], column2index["sepal width"])
4 plt.subplot(1,3,2)
5 plot_cluster(h_clusters, column2index["petal length"], column2index["petal width"])
6 plt.subplot(1,3,3)
7 plot_cluster(h_clusters, column2index["sepal length"], column2index["petal length"])
8 plt.show()
```

```
-----
-
NameError                                Traceback (most recent call last)
<ipython-input-36-94660df5af95> in <module>
      1 plt.subplots(figsize=(20, 5))
      2 plt.subplot(1,3,1)
----> 3 plot_cluster(h_clusters, column2index["sepal length"], column2index["sepal width"])
      4 plt.subplot(1,3,2)
      5 plot_cluster(h_clusters, column2index["petal length"], column2index["petal width"])
```

NameError: name 'plot\_cluster' is not defined

