

```
In [57]: # Cluster on operational packages.
import numpy as np
import pandas as pd

# Important tools for modeling and evaluation.
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_score
from sklearn.preprocessing import StandardScaler

# Import visualization packages.
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [57]: penguins = pd.read_csv('C:\Users\HP\Desktop\Advance Data Analyst\6. The nuts and bolts of ML\3. Module 3\2. Evaluate K Means\Files\penguins.csv')
penguins.head(10)
```

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	Adelie	Torgersen	39.1	18.7	181.0	3750.0	male
1	Adelie	Torgersen	39.5	17.4	186.0	3800.0	female
2	Adelie	Torgersen	40.3	18.0	195.0	3250.0	female
3	Adelie	Torgersen	NaN	NaN	NaN	NaN	NaN
4	Adelie	Torgersen	36.7	19.3	193.0	3450.0	female
5	Adelie	Torgersen	39.3	20.6	190.0	3650.0	male
6	Adelie	Torgersen	38.9	17.8	181.0	3625.0	female
7	Adelie	Torgersen	39.2	19.6	195.0	4675.0	male
8	Adelie	Torgersen	34.1	18.1	193.0	3475.0	NaN
9	Adelie	Torgersen	42.0	20.2	190.0	4250.0	NaN

```
In [57]: # Review the first 10 rows.
penguins.head(n = 10)
```

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	Adelie	Torgersen	39.1	18.7	181.0	3750.0	male
1	Adelie	Torgersen	39.5	17.4	186.0	3800.0	female
2	Adelie	Torgersen	40.3	18.0	195.0	3250.0	female
3	Adelie	Torgersen	NaN	NaN	NaN	NaN	NaN
4	Adelie	Torgersen	36.7	19.3	193.0	3450.0	female
5	Adelie	Torgersen	39.3	20.6	190.0	3650.0	male
6	Adelie	Torgersen	38.9	17.8	181.0	3625.0	female
7	Adelie	Torgersen	39.2	19.6	195.0	4675.0	male
8	Adelie	Torgersen	34.1	18.1	193.0	3475.0	NaN
9	Adelie	Torgersen	42.0	20.2	190.0	4250.0	NaN

```
In [58]: # Find out how many penguin types there are.
penguins['species'].unique()
```

```
Out[58]: array(['Adelie', 'Chinstrap', 'Gentoo'], dtype=object)
```

```
In [58]: # Find the count of each species type.
penguins['species'].value_counts(dropna = False)
```

species	count
Adelie	152
Gentoo	124
Chinstrap	68
Name: count, dtype: int64	

```
In [58]: # Check for missing values.
penguins.isnull().sum()
```

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	Adelie	Torgersen	39.1	18.7	181.0	3750.0	male
1	Adelie	Torgersen	39.5	17.4	186.0	3800.0	female
2	Adelie	Torgersen	40.3	18.0	195.0	3250.0	female
3	Adelie	Torgersen	NaN	NaN	NaN	NaN	NaN
4	Adelie	Torgersen	36.7	19.3	193.0	3450.0	female
5	Adelie	Torgersen	39.3	20.6	190.0	3650.0	male
6	Adelie	Torgersen	38.9	17.8	181.0	3625.0	female
7	Adelie	Torgersen	39.2	19.6	195.0	4675.0	male
8	Adelie	Torgersen	34.1	18.1	193.0	3475.0	NaN
9	Adelie	Torgersen	42.0	20.2	190.0	4250.0	NaN

```
In [58]: # Drop rows with missing values.
# Save DataFrame in variable 'penguins_subset'.
penguins_subset = penguins.dropna(axis=0).reset_index(drop = True)
```

```
In [58]: # Check for missing values.
penguins_subset.isna().sum()
```

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex
0	Adelie	Torgersen	39.1	18.7	181.0	3750.0	male
1	Adelie	Torgersen	39.5	17.4	186.0	3800.0	female
2	Adelie	Torgersen	40.3	18.0	195.0	3250.0	female
3	Adelie	Torgersen	36.7	19.3	193.0	3450.0	female
4	Adelie	Torgersen	39.3	20.6	190.0	3650.0	male
5	Adelie	Torgersen	38.9	17.8	181.0	3625.0	female
6	Adelie	Torgersen	39.2	19.6	195.0	4675.0	male
7	Adelie	Torgersen	41.1	17.6	182.0	3200.0	female
8	Adelie	Torgersen	38.6	21.2	191.0	3800.0	male
9	Adelie	Torgersen	34.6	21.1	198.0	4400.0	male

```
In [58]: penguins_subset['sex'] = penguins_subset['sex'].str.upper()
```

```
In [58]: # Convert 'sex' column from categorical to numeric.
penguins_subset = pd.get_dummies(penguins_subset, drop_first = True, columns=['sex'])
```

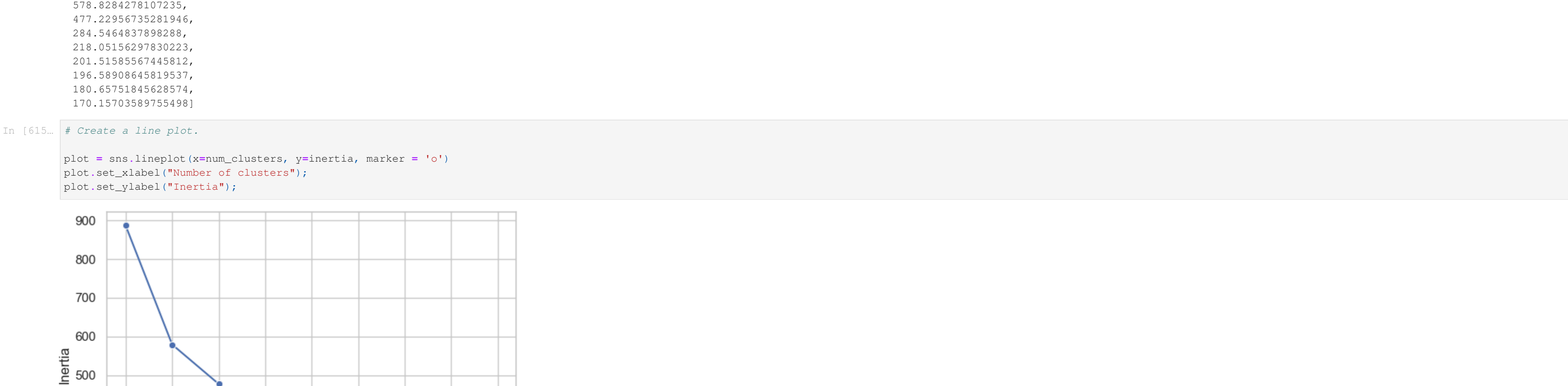
```
In [58]: # Drop the 'island' column.
penguins_subset = penguins_subset.drop(['island'], axis=1)
```

```
In [58]: # Exclude 'species' variable from X
X = penguins_subset.drop(['species'], axis=1)
```

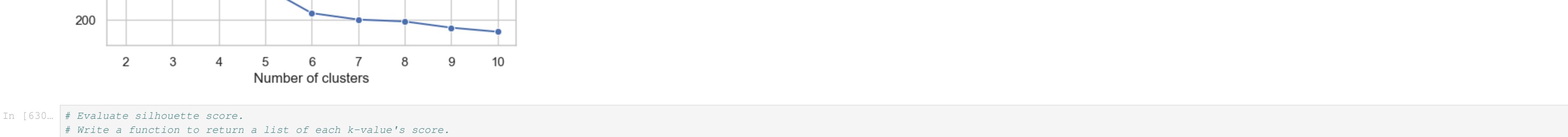
```
In [60]: #Scale the Features.
#Assign the scaled data to variable 'X_scaled'.
X_scaled = StandardScaler().fit_transform(X)
```

```
In [60]: # Fit K-means and evaluate inertia for different values of k.
num_clusters = [i for i in range(2, 11)]
def kmeans_inertia(num_clusters, x_vals):
    """
    Accepts as arguments list of ints and data array.
    Fits a KMeans model where k = each value in the list of ints.
    Returns each k-value's inertia appended to a list.
    """
    inertia = []
    for num in num_clusters:
        kms = KMeans(n_clusters=num, random_state=42)
        kms.fit(X_vals)
        inertia.append(kms.inertia_)
    return inertia
```

```
In [60]: # Return a list of inertia for k=2 to 10.
inertia = kmeans_inertia(num_clusters, X_scaled)
inertia
```

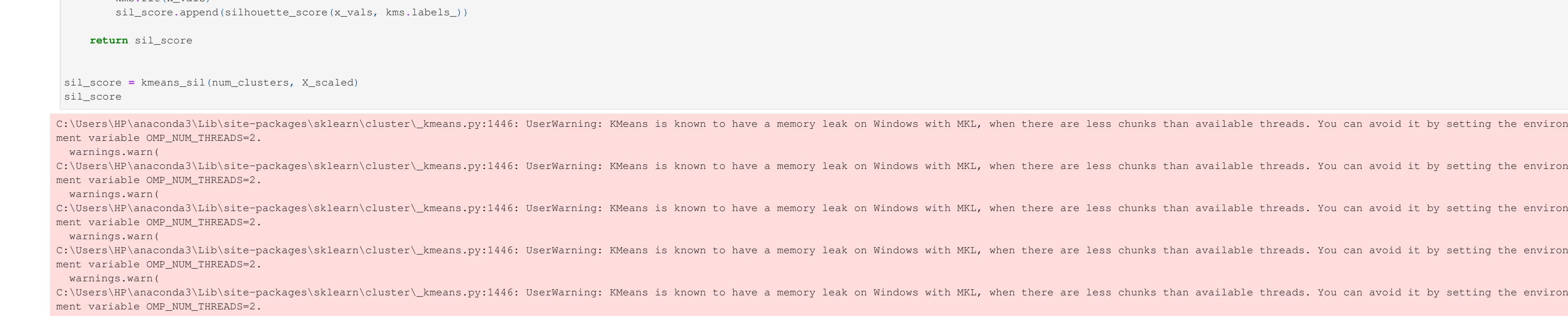


```
In [60]: # Create a line plot.
plot = sns.lineplot(x=num_clusters, y=inertia, marker = 'o')
plt.set_xlabel('Number of clusters')
plt.set_ylabel('Inertia')
```

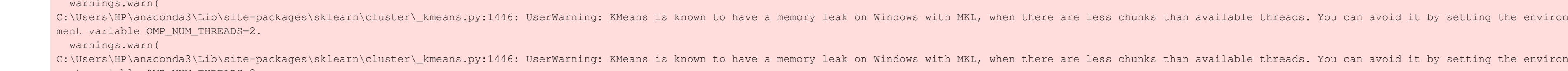


```
In [60]: # Evaluate silhouette score.
# Write a function to return a list of each k-value's score.
def kmeans_sil(num_clusters, x_vals):
    """
    Accepts as arguments list of ints and data array.
    Fits a KMeans model where k = each value in the list of ints.
    Calculates a silhouette score for each k value.
    Returns each k-value's silhouette score appended to a list.
    """
    sil_score = []
    for num in num_clusters:
        kms = KMeans(n_clusters=num, random_state=42)
        kms.fit(X_vals)
        sil_score.append(silhouette_score(x_vals, kms.labels_))
    return sil_score
```

```
sil_score = kmeans_sil(num_clusters, X_scaled)
sil_score
```



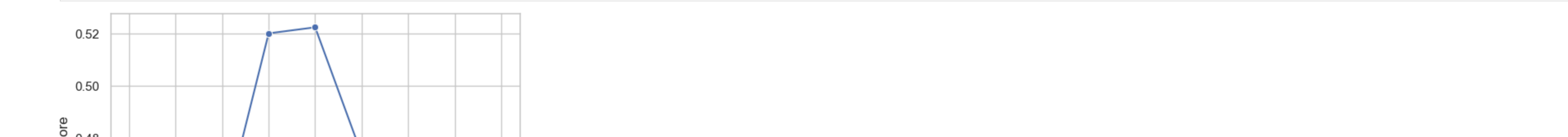
```
In [60]: # Create a line plot.
plot = sns.lineplot(x=num_clusters, y=sil_score, marker = 'o')
plt.set_xlabel('# of clusters')
plt.set_ylabel('Silhouette Score')
```



```
In [60]: # Fit a 6-cluster model.
kmeans6 = KMeans(n_clusters=6, random_state=42)
kmeans6.fit(X_scaled)
```

```
Out[60]: 0.4439808835305243,
0.43101024097188366,
0.448989212001027,
0.51998574860868,
0.5220860608437773,
0.4738633056429317,
0.471544346463867,
0.4180514949494056,
0.418306343951951
```

```
In [60]: # Create a line plot.
plot = sns.lineplot(x=num_clusters, y=sil_score, marker = 'o')
plt.set_xlabel('# of clusters')
plt.set_ylabel('Silhouette Score')
```



```
In [60]: # Verify if any 'cluster' can be differentiated by 'species'.
penguins_subset.groupby(by=['cluster', 'species'])['size']
```

cluster	species	size
0	Chinstrap	32
1	Gentoo	58
2	Adelie	73
3	Chinstrap	2
4	Adelie	2
5	Chinstrap	34
6	Gentoo	61
7	Adelie	71

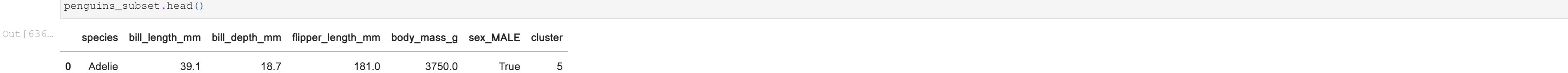
```
In [60]: # Verify if any 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().sort_values(ascending = False)
```

cluster	species	sex_MALE	size
0	Adelie	False	73
1	Adelie	True	71
2	Gentoo	True	61
3	Gentoo	False	58
4	Chinstrap	True	34
5	Chinstrap	False	32
6	Adelie	True	2

```
In [64]: # Print unique labels.
print('Unique labels:', np.unique(kmeans6.labels_))
Unique labels: [0 1 2 3 4 5]
```

```
In [63]: # Create a new column 'cluster'.
penguins_subset['cluster'] = kmeans6.labels_
penguins_subset.head()
```

	species	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g	sex_MALE	cluster
0	Adelie	39.1	18.7	181.0	3750.0	True	5
1	Adelie	39.5	17.4	186.0	3800.0	False	2
2	Adelie	40.3	18.0	195.0	3250.0	False	2
3	Adelie	36.7	19.3	193.0	3450.0	False	2
4	Adelie	39.3	20.6	190.0	3650.0	True	5



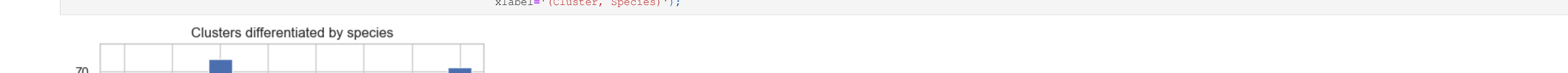
```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().sort_values(ascending = False)
```

cluster	species	sex_MALE	size
0	Adelie	False	73
1	Adelie	True	71
2	Gentoo	True	61
3	Gentoo	False	58
4	Chinstrap	True	34
5	Chinstrap	False	32
6	Adelie	True	2

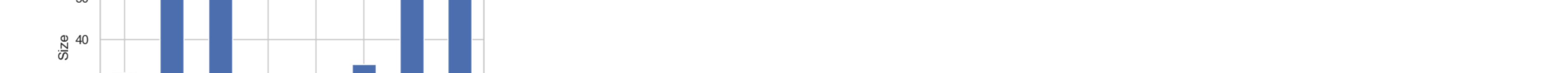
```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



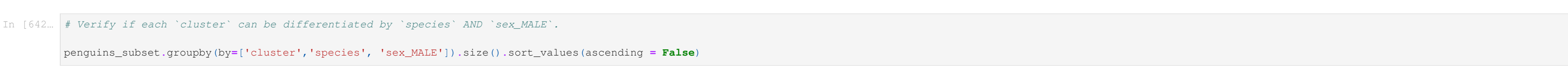
```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



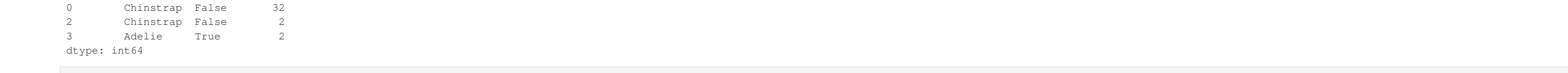
```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



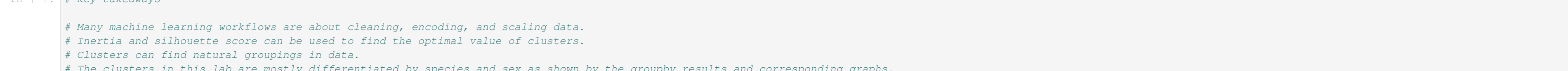
```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



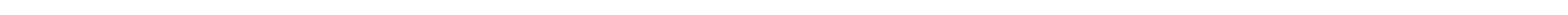
```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



```
In [64]: # Verify if each 'cluster' can be differentiated by 'species' AND 'sex_MALE'.
penguins_subset.groupby(by=['cluster', 'species', 'sex_MALE']).size().unstack(level = 'species', fill_value=0).plot.bar(title='Clusters differentiated by species and sex',
figsize=(6, 5),
ylabel='Size',
xlabel='(Cluster, Sex)')
```



# What summary would you provide to stakeholders?

# The K-means clustering enabled this data to be effectively grouped. It helped identify patterns that can educate team members about penguins.

# The success of the cluster results suggests that the organization can apply clustering to other projects and continue augmenting employee education.