PROJECT: CLUSTERING ANTARCTIC PENGUIN SPECIES



source: @allison_horst https://github.com/allisonhorst/penguins

You have been asked to support a team of researchers who have been collecting data about penguins in Antartica! The data is available in csv-Format as `penguins.csv`

**Origin of this data** : Data were collected and made available by Dr. Kristen Gorman and the Palmer Station, Antarctica LTER, a member of the Long Term Ecological Research Network.

**The dataset consists of 5 columns.**

| Column | Description |
| --- | --- |
| culmen_length_mm | culmen length (mm) |
| culmen_depth_mm | culmen depth (mm) |
| flipper_length_mm | flipper length (mm) |
| body_mass_g | body mass (g) |
| sex | penguin sex |

Unfortunately, they have not been able to record the species of penguin, but they know that there are **at least three** species that are native to the region: **Adelie**, **Chinstrap**, and **Gentoo.** Your task is to apply your data science skills to help them identify groups in the dataset!

```python
# Import Required Packages
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

# Loading and examining the dataset
penguins_df = pd.read_csv("penguins.csv")
penguins_df.head()
penguins_df.info()

penguins_df = pd.get_dummies(penguins_df, dtype='int')

scaler = StandardScaler()
X = scaler.fit_transform(penguins_df)
penguins_preprocessed = pd.DataFrame(data=X, columns=penguins_df.columns)
penguins_preprocessed.head(10)

inertia = []
for k in range(1, 10):
    kmeans = KMeans(n_clusters=k, random_state=42).fit(penguins_preprocessed)
    inertia.append(kmeans.inertia_)
plt.plot(range(1, 10), inertia, marker='o')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.title('Elbow Method')
plt.show()
n_clusters = 4

kmeans = KMeans(n_clusters=n_clusters, random_state=42).fit(penguins_preprocessed)
penguins_df['label'] = kmeans.labels_
plt.scatter(penguins_df['label'], penguins_df['culmen_length_mm'], c=kmeans.labels_,
cmap='viridis')
plt.xlabel('Cluster')
plt.ylabel('culmen_length_mm')
plt.xticks(range(int(penguins_df['label'].min()), int(penguins_df['label'].max()) +
1))
plt.title(f'K-means Clustering (K={n_clusters})')
plt.show()

numeric_columns = ['culmen_length_mm', 'culmen_depth_mm', 'flipper_length_mm',
'label']
stat_penguins = penguins_df[numeric_columns].groupby('label').mean()
stat_penguins

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 332 entries, 0 to 331
```
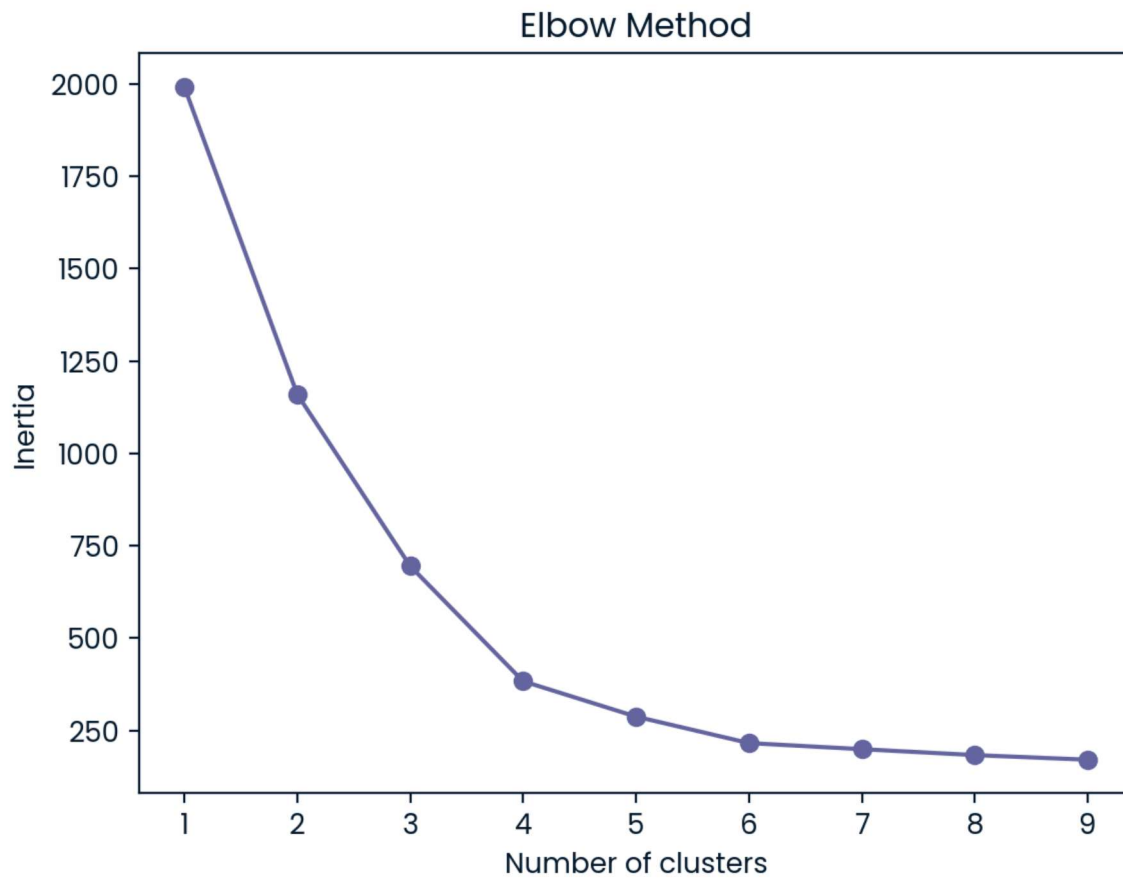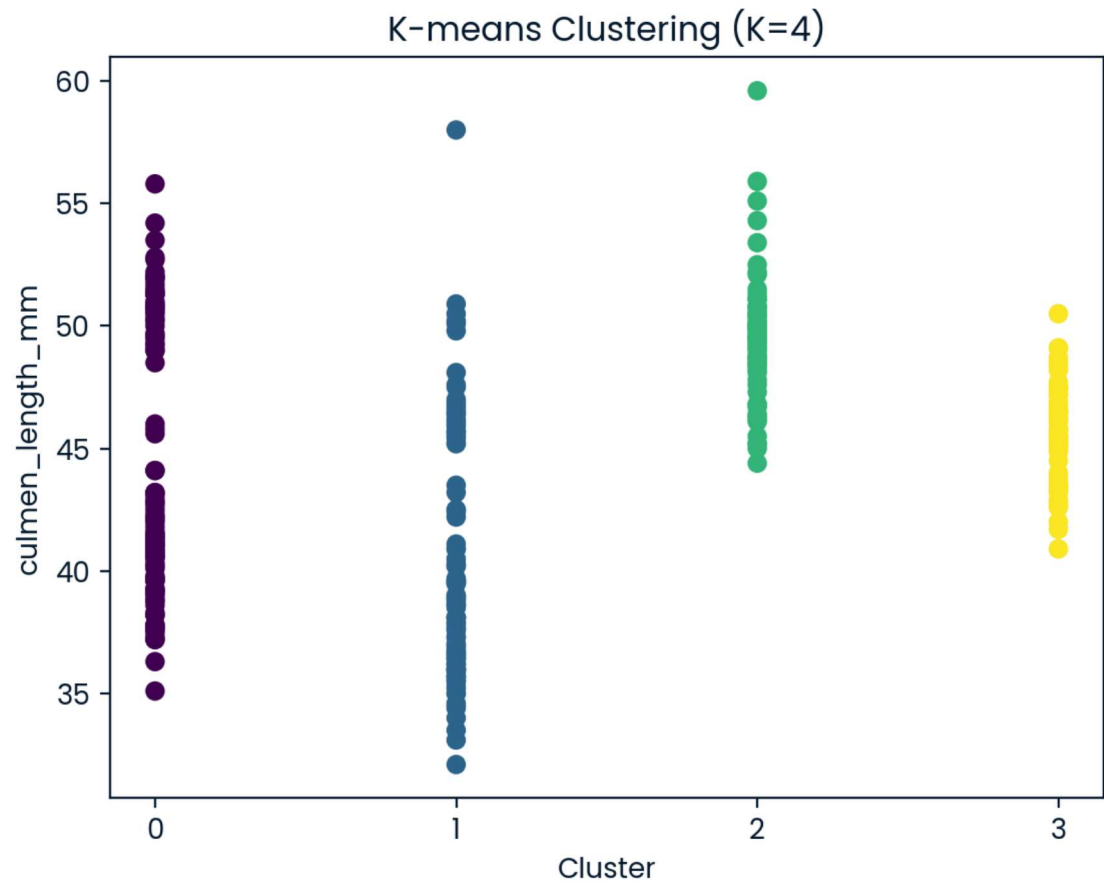
```
Data columns (total 5 columns):
 #   Column             Non-Null Count   Dtype
---  ------             --------------   -----
 0   culmen_length_mm   332 non-null     float64
 1   culmen_depth_mm    332 non-null     float64
 2   flipper_length_mm  332 non-null     float64
 3   body_mass_g        332 non-null     float64
 4   sex                332 non-null     object
dtypes: float64(4), object(1)
memory usage: 13.1+ KB
```



Elbow Method

## K-means Clustering (K=4)



| | culmen_len... | | culmen_d... | | flipper_length... | |
|---|---|---|---|---|---|---|
| 0 | 43.8783018868 | | 19.1113207547 | | 194.7641509434 | |
| 1 | 40.2177570093 | | 17.6112149533 | | 189.046728972 | |
| 2 | 49.4737704918 | | 15.7180327869 | | 221.5409836066 | |
| 3 | 45.5637931034 | | 14.2379310345 | | 212.7068965517 | |

Rows: 4                                                                    ⤢ Expand