

AI and ML task 8

objective

1.Load and visualize dataset (optional PCA for 2D view). 2.Fit K-Means and assign cluster labels. 3.Use the Elbow Method to find optimal K. 4.Visualize clusters with color-coding. 5.Evaluate clustering using Silhouette Score.

Importing libraries

```
In [5]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

In [7]: data=pd.read_csv('Mall_Customers.csv')

In [9]: data

Out[9]:
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)	
	0	1	Male	19	15	39
	1	2	Male	21	15	81
	2	3	Female	20	16	6
	3	4	Female	23	16	77
	4	5	Female	31	17	40
...
	195	196	Female	35	120	79
	196	197	Female	45	126	28
	197	198	Male	32	126	74
	198	199	Male	32	137	18
	199	200	Male	30	137	83

200 rows × 5 columns

```
In [11]: data.shape

Out[11]: (200, 5)

In [13]: data.describe(include='all')

Out[13]:
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200	200.000000	200.000000	200.000000
unique	NaN	2	NaN	NaN	NaN
top	NaN	Female	NaN	NaN	NaN
freq	NaN	112	NaN	NaN	NaN
mean	100.500000	NaN	38.850000	60.560000	50.200000
std	57.879185	NaN	13.969007	26.264721	25.823522
min	1.000000	NaN	18.000000	15.000000	1.000000
25%	50.750000	NaN	28.750000	41.500000	34.750000
50%	100.500000	NaN	36.000000	61.500000	50.000000
75%	150.250000	NaN	49.000000	78.000000	73.000000
max	200.000000	NaN	70.000000	137.000000	99.000000

```
In [15]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column              Non-Null Count  Dtype
---  ---
0   CustomerID          200 non-null    int64
1   Gender              200 non-null    object
2   Age                 200 non-null    int64
3   Annual Income (k$)  200 non-null    int64
4   Spending Score (1-100) 200 non-null    int64
dtypes: int64(4), object(1)
memory usage: 7.9+ KB

In [17]: data.isnull()

Out[17]:
```

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	False	False	False	False	False
1	False	False	False	False	False
2	False	False	False	False	False
3	False	False	False	False	False
4	False	False	False	False	False
...
195	False	False	False	False	False
196	False	False	False	False	False
197	False	False	False	False	False
198	False	False	False	False	False
199	False	False	False	False	False

200 rows × 5 columns

```
In [19]: data.isnull().sum()

Out[19]: CustomerID    0
Gender              0
Age                 0
Annual Income (k$)  0
Spending Score (1-100) 0
dtype: int64
```

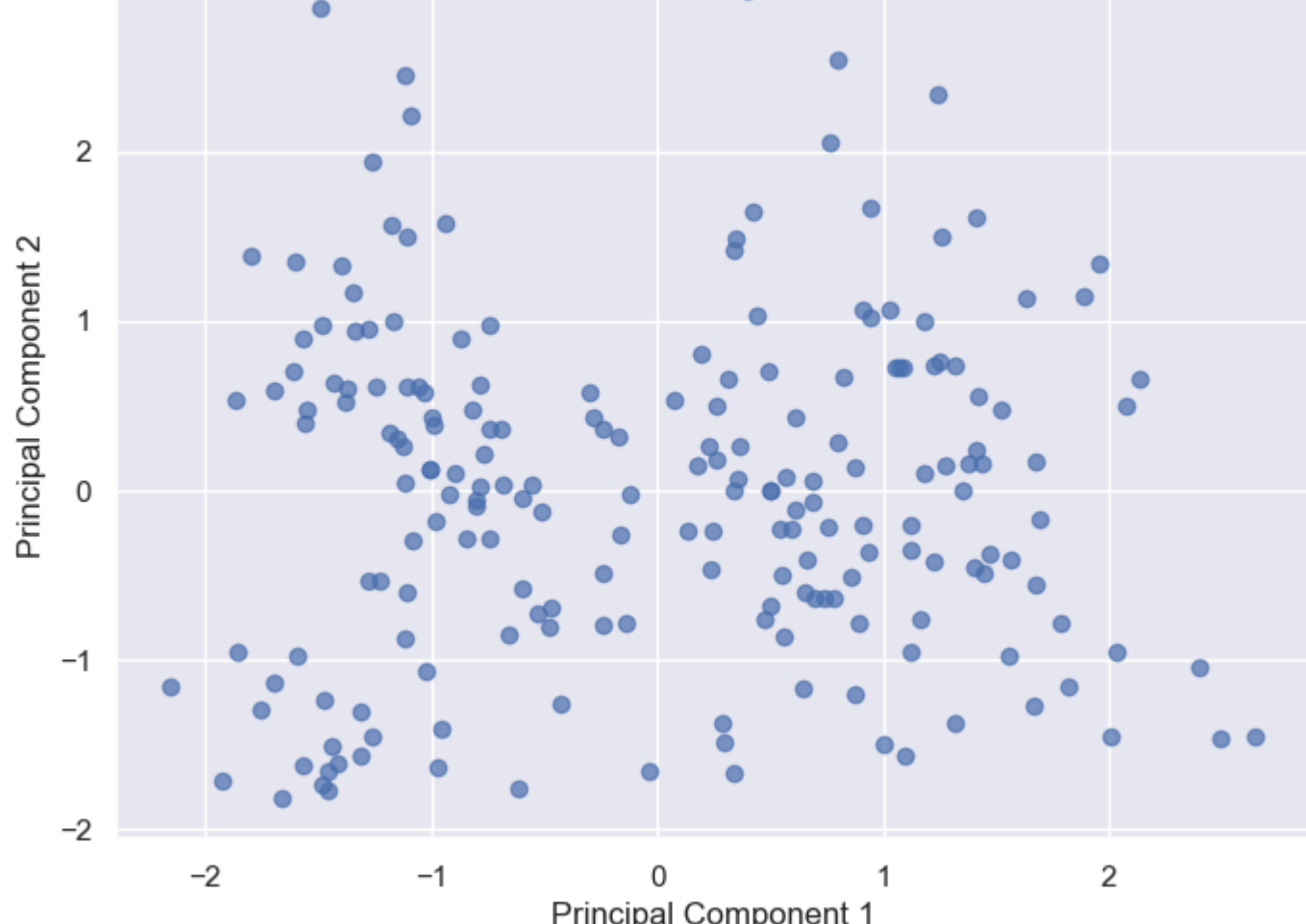
Load and visualize dataset (optional PCA for 2D view).

```
In [24]: features = ['Age', 'Annual Income (k$)', 'Spending Score (1-100)']
X = data[features]

In [26]: # Standardize the data
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

In [28]: # Apply PCA to reduce to 2 dimensions for visualization
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X_scaled)

In [30]: plt.figure(figsize=(8, 6))
plt.scatter(X_pca[:, 0], X_pca[:, 1], alpha=0.7)
plt.title('PCA Projection of Mall Customers')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.grid(True)
plt.show()
```



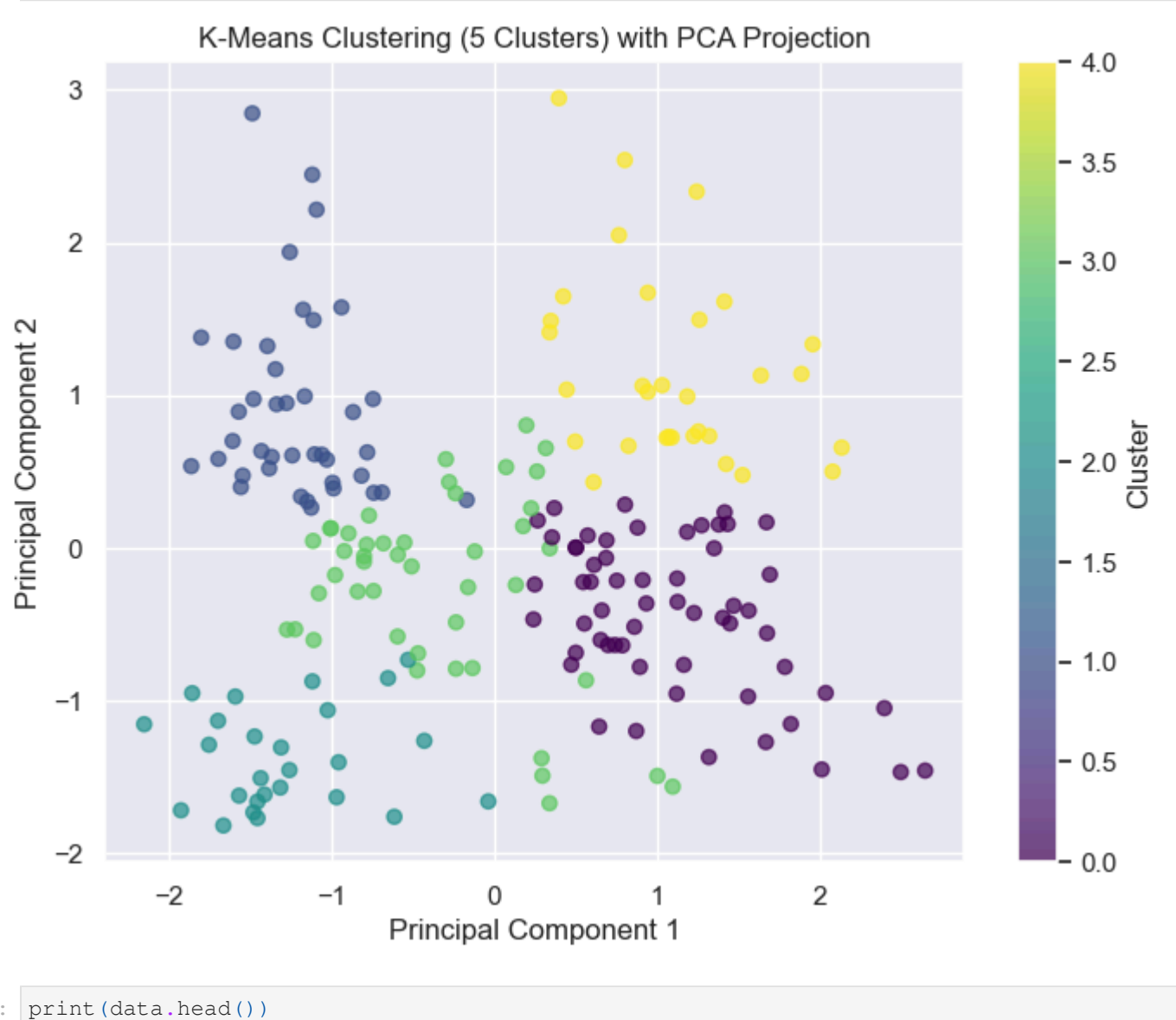
Fit K-Means and assign cluster labels.

```
In [46]: from sklearn.cluster import KMeans
import warnings
warnings.filterwarnings('ignore')

In [48]: # Fit K-Means and predict clusters
kmeans = KMeans(n_clusters=5, random_state=42)
clusters = kmeans.fit_predict(X_scaled)

In [50]: data['Cluster'] = clusters

In [52]: plt.figure(figsize=(8, 6))
scatter = plt.scatter(X_pca[:, 0], X_pca[:, 1], c=clusters, cmap='viridis', alpha=0.7)
plt.title('K-Means Clustering (5 Clusters) with PCA Projection')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.grid(True)
plt.colorbar(scatter, label='Cluster')
plt.show()
```



```
In [54]: print(data.head())

CustomerID  Gender  Age  Annual Income (k$)  Spending Score (1-100) \
0           1    Male   19                15                  39
1           2    Male   21                15                  81
2           3    Female  20                16                   6
3           4    Female  23                16                 77
4           5    Female  31                17                 40

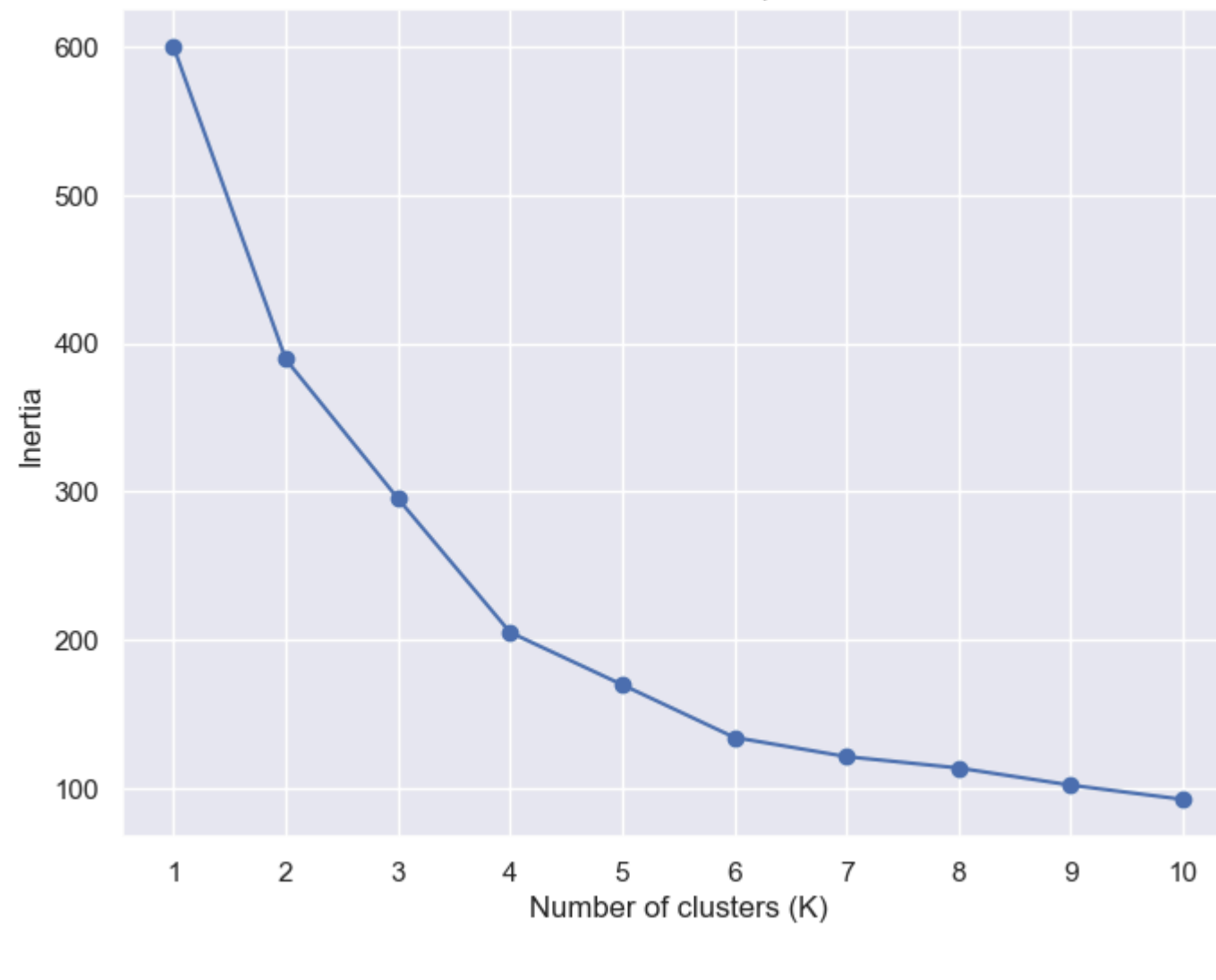
Cluster
0           2
1           2
2           3
3           2
4           2
```

Use the Elbow Method to find optimal K.

```
In [56]: inertia = []
K_range = range(1, 11)

for k in K_range:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(X_scaled)
    inertia.append(kmeans.inertia_)

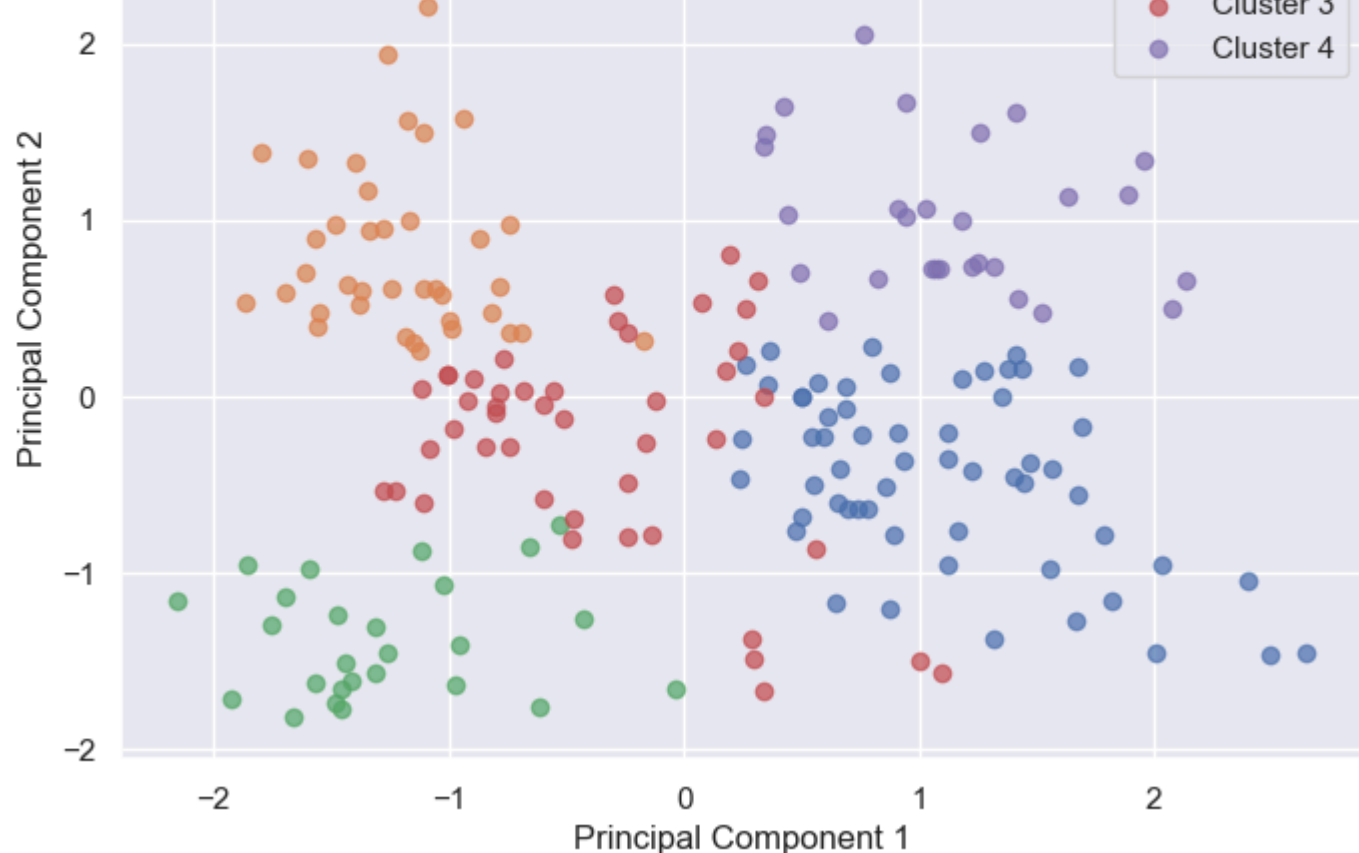
In [62]: # Plot the Elbow curve
plt.figure(figsize=(8, 6))
plt.plot(K_range, inertia, marker='o')
plt.title('Elbow Method For Optimal K')
plt.xlabel('Number of clusters (K)')
plt.ylabel('Inertia')
plt.grid(True)
plt.xticks(K_range)
plt.show()
```



Visualize clusters with color-coding.

```
In [66]: plt.figure(figsize=(8, 6))
for cluster in range(5):
    plt.scatter(
        X_pca[clusters == cluster, 0],
        X_pca[clusters == cluster, 1],
        label=f'Cluster {cluster}',
        alpha=0.7
    )

plt.title('Mall Customers Cluster Visualization (K=5)')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.legend()
plt.grid(True)
plt.show()
```



Evaluate clustering using Silhouette Score.

```
In [71]: from sklearn.metrics import silhouette_score

In [73]: sil_score = silhouette_score(X_scaled, clusters)

sil_score
```

