

```
import numpy as np
import pandas as pd
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
```

1. PCA

```
file_path = 'winequality-white.csv'
df = pd.read_csv(file_path, sep=';')
```

```
df.head()
```

```
{
  "summary": {
    "name": "df",
    "rows": 4898,
    "fields": [
      {
        "column": "fixed acidity",
        "properties": {
          "dtype": "number",
          "std": 0.8438682276875188,
          "min": 3.8,
          "max": 14.2,
          "num_unique_values": 68,
          "samples": [
            10.3,
            5.8,
            6.2
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "volatile acidity",
        "properties": {
          "dtype": "number",
          "std": 0.10079454842486428,
          "min": 0.08,
          "max": 1.1,
          "num_unique_values": 125,
          "samples": [
            0.14,
            0.595,
            0.13
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "citric acid",
        "properties": {
          "dtype": "number",
          "std": 0.12101980420298301,
          "min": 0.0,
          "max": 1.66,
          "num_unique_values": 87,
          "samples": [
            0.64,
            0.36,
            0.24
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "residual sugar",
        "properties": {
          "dtype": "number",
          "std": 5.072057784014864,
          "min": 0.6,
          "max": 65.8,
          "num_unique_values": 310,
          "samples": [
            15.5,
            19.25,
            3.3
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "chlorides",
        "properties": {
          "dtype": "number",
          "std": 0.02184796809372882,
          "min": 0.009,
          "max": 0.346,
          "num_unique_values": 160,
          "samples": [
            0.133,
            0.015
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "free sulfur dioxide",
        "properties": {
          "dtype": "number",
          "std": 17.007137325232566,
          "min": 2.0,
          "max": 289.0,
          "num_unique_values": 132,
          "samples": [
            24.0,
            122.5,
            7.0
          ],
          "semantic_type": "",
          "description": ""
        },
        "column": "total sulfur dioxide",
        "properties": {
          "dtype": "number",
          "std": 0.02184796809372882,
          "min": 0.009,
          "max": 0.346,
          "num_unique_values": 160,
          "samples": [
            0.133,
            0.015
          ],
          "semantic_type": "",
          "description": ""
        }
      ]
    }
  }
}
```

```

{"number": 9.0, "std": 42.49806455414294, "min": 9.0, "max": 440.0, "num_unique_values": 251, "samples": [260.0, 63.0, 70.0], "semantic_type": "", "description": ""}, {"column": "density", "properties": {"dtype": "number", "std": 0.0029909069169369354, "min": 0.98711, "max": 1.03898, "num_unique_values": 890, "samples": [0.99362, 0.99388, 0.9929]}, "semantic_type": "", "description": ""}, {"column": "pH", "properties": {"dtype": "number", "std": 0.1510005996150667, "min": 2.72, "max": 3.82, "num_unique_values": 103, "samples": [3.34, 3.41, 3.49]}, "semantic_type": "", "description": ""}, {"column": "sulphates", "properties": {"dtype": "number", "std": 0.11412583394883138, "min": 0.22, "max": 1.08, "num_unique_values": 79, "samples": [0.41, 0.45, 0.46]}, "semantic_type": "", "description": ""}, {"column": "alcohol", "properties": {"dtype": "number", "std": 1.2306205677573183, "min": 8.0, "max": 14.2, "num_unique_values": 103, "samples": [12.6, 11.366666666666667, 10.033333333333333]}, "semantic_type": "", "description": ""}, {"column": "quality", "properties": {"dtype": "number", "std": 0, "min": 3, "max": 9, "num_unique_values": 7, "samples": [6, 5, 3]}, "semantic_type": "", "description": ""}]
{"type": "dataframe", "variable_name": "df"}

```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 4898 entries, 0 to 4897
```

```
Data columns (total 12 columns):
```

#	Column	Non-Null Count	Dtype
0	fixed acidity	4898 non-null	float64
1	volatile acidity	4898 non-null	float64
2	citric acid	4898 non-null	float64
3	residual sugar	4898 non-null	float64
4	chlorides	4898 non-null	float64
5	free sulfur dioxide	4898 non-null	float64
6	total sulfur dioxide	4898 non-null	float64
7	density	4898 non-null	float64
8	pH	4898 non-null	float64

```
9   sulphates          4898 non-null   float64
10  alcohol            4898 non-null   float64
11  quality            4898 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 459.3 KB

features = df.drop(columns=['quality'])

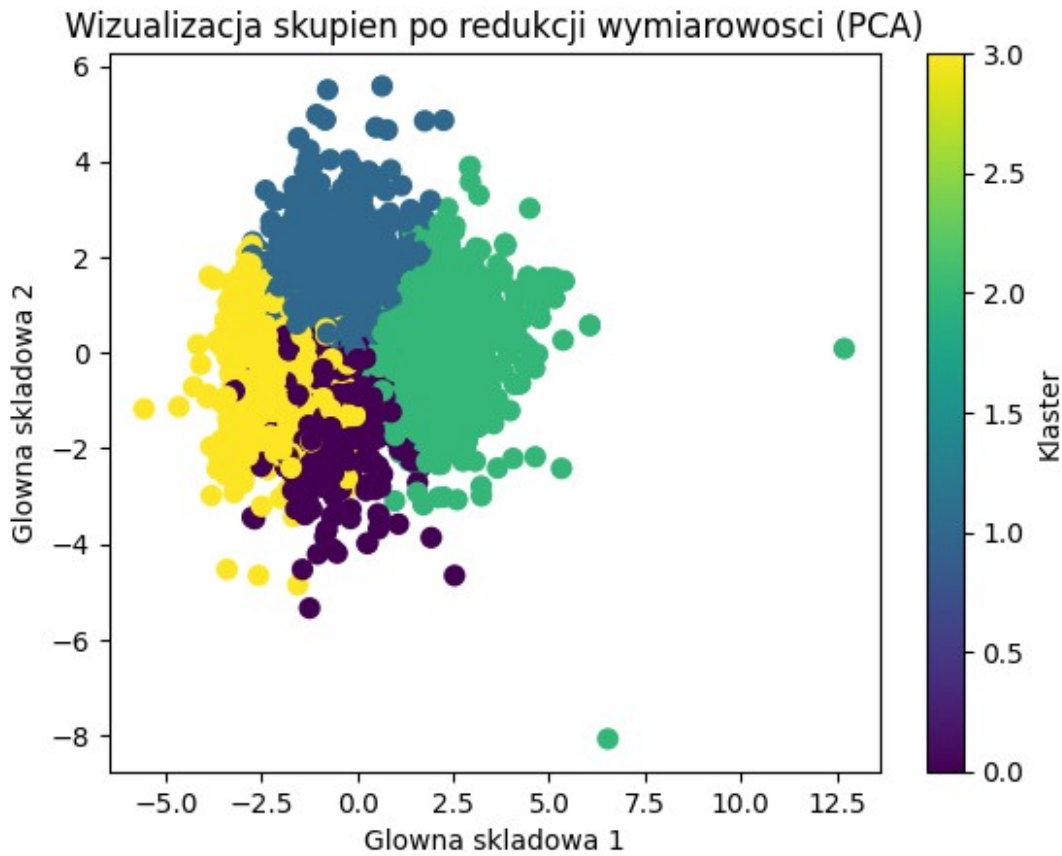
scaler = StandardScaler()
features_scaled = scaler.fit_transform(features)

pca = PCA(n_components=2)

features_pca = pca.fit_transform(features_scaled)

kmeans_pca = KMeans(n_clusters=4, random_state=42)
labels_pca = kmeans_pca.fit_predict(features_pca)

plt.scatter(features_pca[:, 0], features_pca[:, 1], c=labels_pca,
            cmap='viridis', s=50)
plt.xlabel('Główna składowa 1')
plt.ylabel('Główna składowa 2')
plt.title('Wizualizacja skupień po redukcji wymiarowości (PCA)')
plt.colorbar(label='Klaster')
plt.show()
```



#PORÓWNANIE WYNIKÓW PRZED I PO REDUKCJI WIELOWYMIAROWOŚCI - KORZYSTAMY ZE WSPÓŁCZYNNIKA SILHOUETTE

```
from sklearn.metrics import silhouette_score
```

```
silhouette_pca = silhouette_score(features_pca, labels_pca)
print(f"Silhouette Score (PCA-transformed data): {silhouette_pca:.2f}")
```

Silhouette Score (PCA-transformed data): 0.38

```
kmeans_original = KMeans(n_clusters=3, random_state=42)
labels_original = kmeans_original.fit_predict(features_scaled)
silhouette_original = silhouette_score(features_scaled, labels_original)
```

```
print(f"Silhouette Score (original data): {silhouette_original:.2f}")
```

Silhouette Score (original data): 0.14

2. METODY NIEHIERARCHICZNE

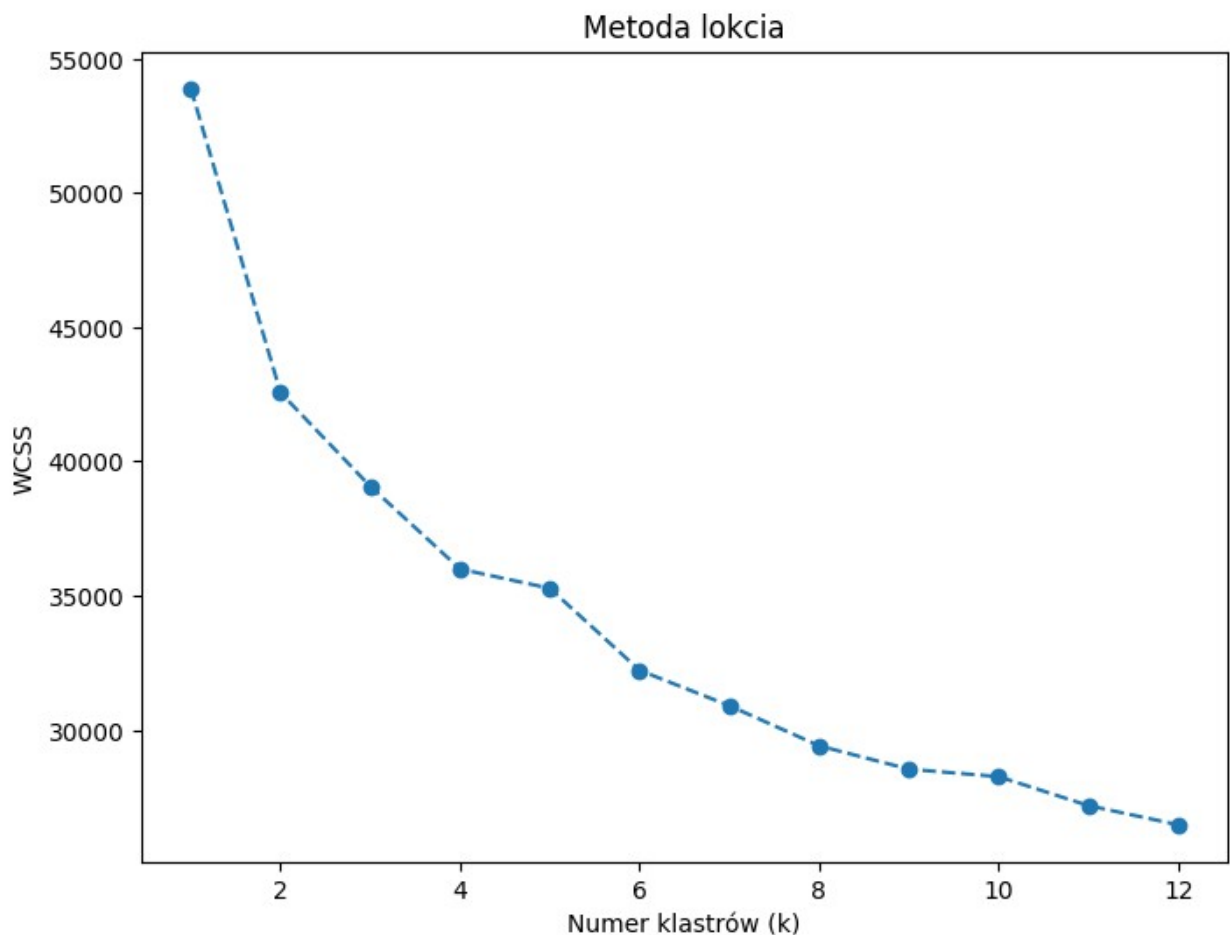
```
wcss = []
k_values = range(1, 13)
for k in k_values:
```

```

kmeans = KMeans(n_clusters=k, random_state=42)
kmeans.fit(features_scaled)
wcss.append(kmeans.inertia_)

plt.figure(figsize=(8, 6))
plt.plot(k_values, wcss, marker='o', linestyle='--')
plt.title('Metoda lokcia')
plt.xlabel('Numer klastrów (k)')
plt.ylabel('WCSS')
plt.show()

```



```

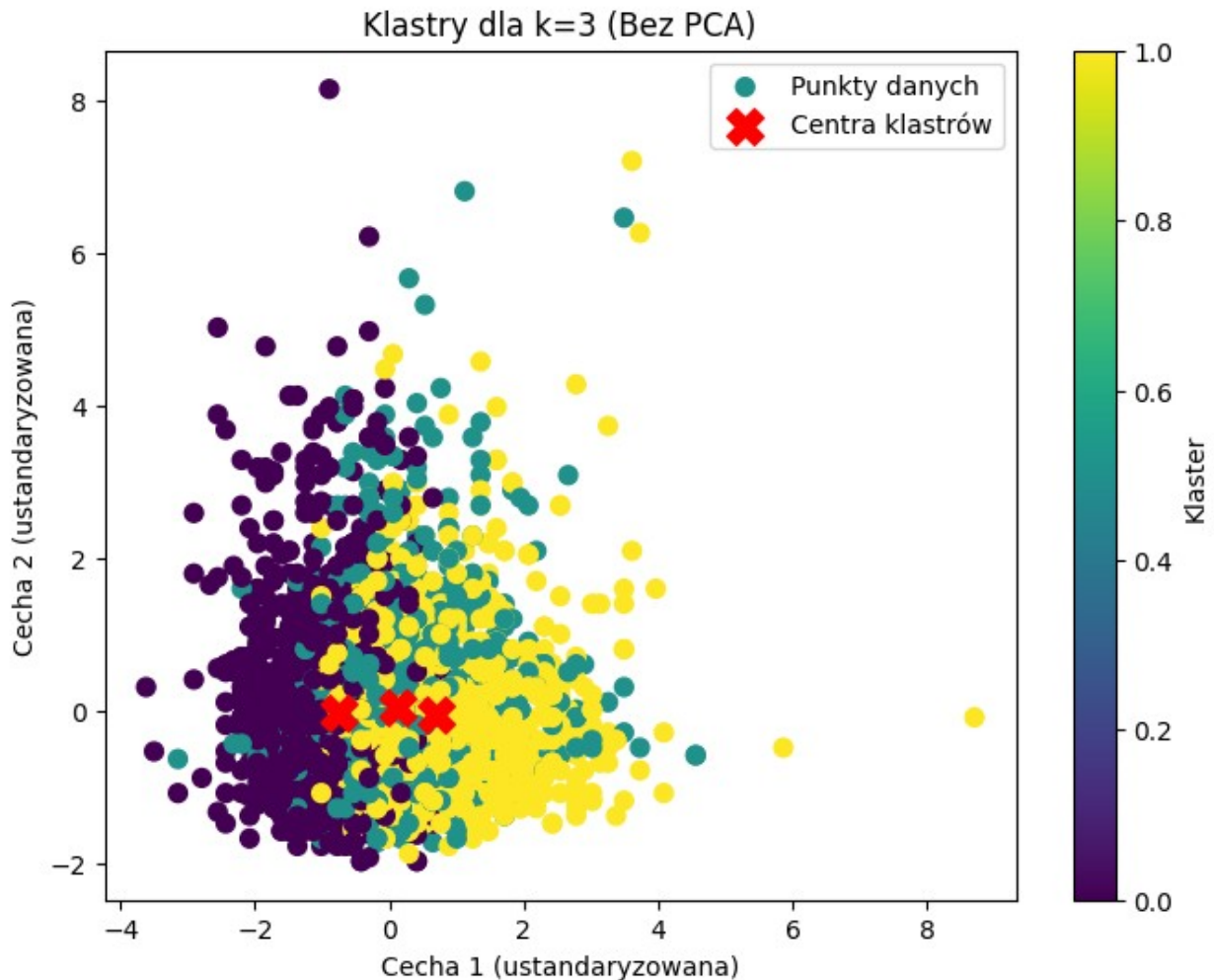
#wybieramy k=3
optimal_k = 3

kmeans_optimal = KMeans(n_clusters=optimal_k, random_state=42)
labels_optimal = kmeans_optimal.fit_predict(features_scaled)

plt.figure(figsize=(8, 6))
plt.scatter(features_scaled[:, 0], features_scaled[:, 1],
            c=labels_optimal, cmap='viridis', s=50, label='Punkty danych')

```

```
plt.scatter(kmeans_optimal.cluster_centers_[0, 0],
            kmeans_optimal.cluster_centers_[0, 1],
            c='red', marker='X', s=200, label='Centra klastrów')
plt.xlabel('Cecha 1 (ustandaryzowana)')
plt.ylabel('Cecha 2 (ustandaryzowana)')
plt.title(f'Klastry dla k={optimal_k} (Bez PCA)')
plt.colorbar(label='Klaster')
plt.legend()
plt.show()
```

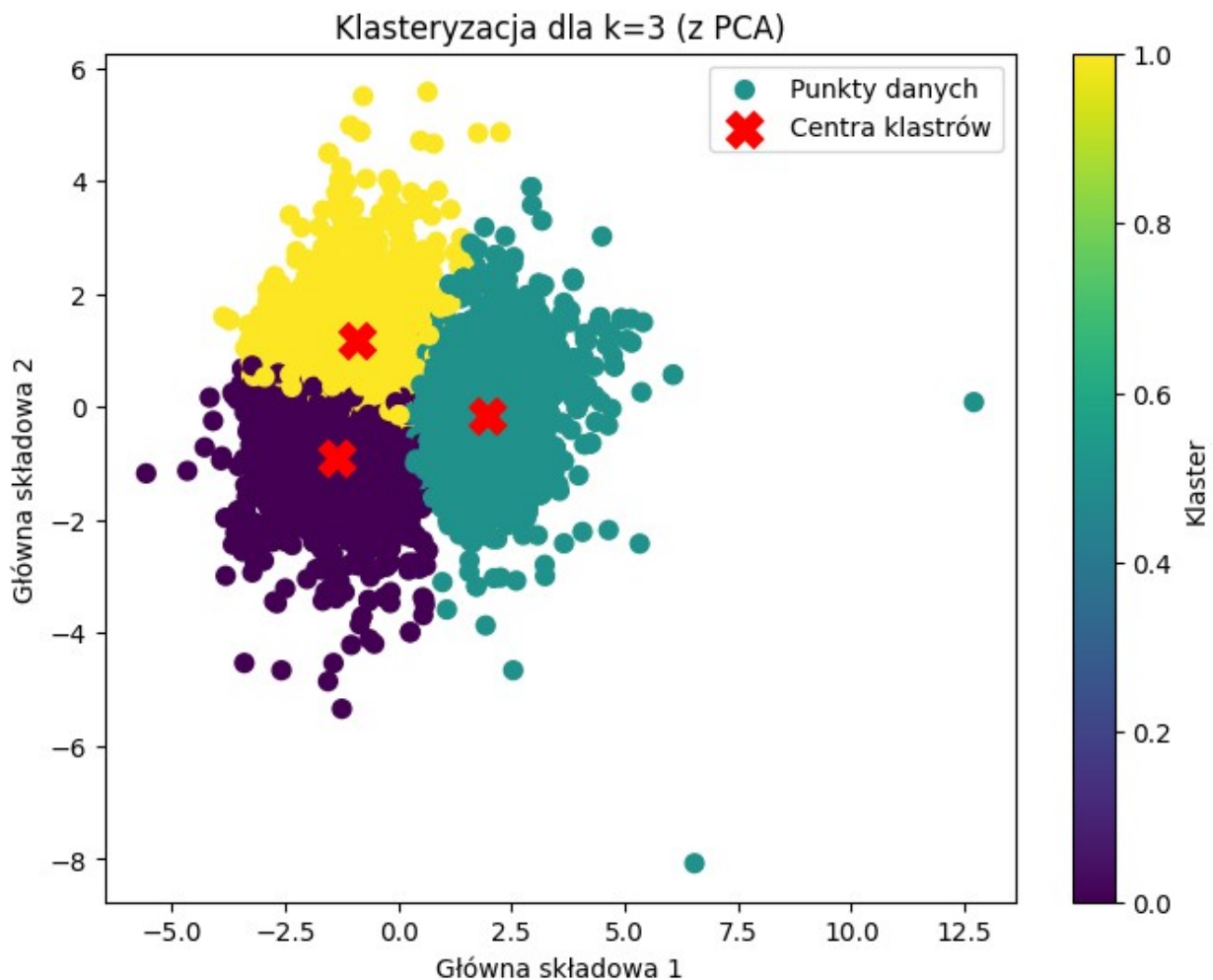


```
kmeans_optimal = KMeans(n_clusters=optimal_k, random_state=42)
labels_optimal = kmeans_optimal.fit_predict(features_scaled)

pca = PCA(n_components=2)
features_pca = pca.fit_transform(features_scaled)

cluster_centers_pca = pca.transform(kmeans_optimal.cluster_centers_)
```

```
plt.figure(figsize=(8, 6))
plt.scatter(features_pca[:, 0], features_pca[:, 1], c=labels_optimal,
            cmap='viridis', s=50, label='Punkty danych')
plt.scatter(cluster_centers_pca[:, 0], cluster_centers_pca[:, 1],
            c='red', marker='X', s=200, label='Centra klastrów')
plt.xlabel('Główna składowa 1')
plt.ylabel('Główna składowa 2')
plt.title(f'Klasteryzacja dla k={optimal_k} (z PCA)')
plt.colorbar(label='Klaster')
plt.legend()
plt.show()
```



```
import scipy.cluster.hierarchy as sch
from scipy.cluster.hierarchy import dendrogram, linkage, fcluster

cechy = df.drop(columns=['quality'])
scaler = StandardScaler()
dane_standaryzowane = scaler.fit_transform(cechy)
```

```

dane_standaryzowane = dane_standaryzowane[:1000]

metody_laczenia = ['ward', 'single', 'complete']

for i, metoda in enumerate(metody_laczenia):
    plt.figure(figsize=(15, 10))
    linked = linkage(dane_standaryzowane, method=metoda)
    dendrogram(linked, orientation='top', distance_sort='descending',
show_leaf_counts=False)
    plt.title(f'Dendrogram (metoda {metoda})')
    plt.xlabel('Próbki')
    plt.ylabel('Odległość')
    plt.tight_layout()
    plt.show()

linked = linkage(dane_standaryzowane, method='ward')
klastry_hierarchiczne = fcluster(linked, t=3, criterion='maxclust')

plt.scatter(dane_standaryzowane[:, 0], dane_standaryzowane[:, 1],
c=klastry_hierarchiczne, cmap='viridis', s=50)
plt.title('Klasteryzacja hierarchiczna (metoda Warda)')
plt.xlabel('Cecha 1 (standaryzowana)')
plt.ylabel('Cecha 2 (standaryzowana)')
plt.show()

kmeans = KMeans(n_clusters=3, random_state=42)
klastry_kmeans = kmeans.fit_predict(dane_standaryzowane)

plt.scatter(dane_standaryzowane[:, 0], dane_standaryzowane[:, 1],
c=klastry_kmeans, cmap='viridis', s=50)
plt.title('Klasteryzacja K-Means')
plt.xlabel('Cecha 1 (standaryzowana)')
plt.ylabel('Cecha 2 (standaryzowana)')
plt.show()

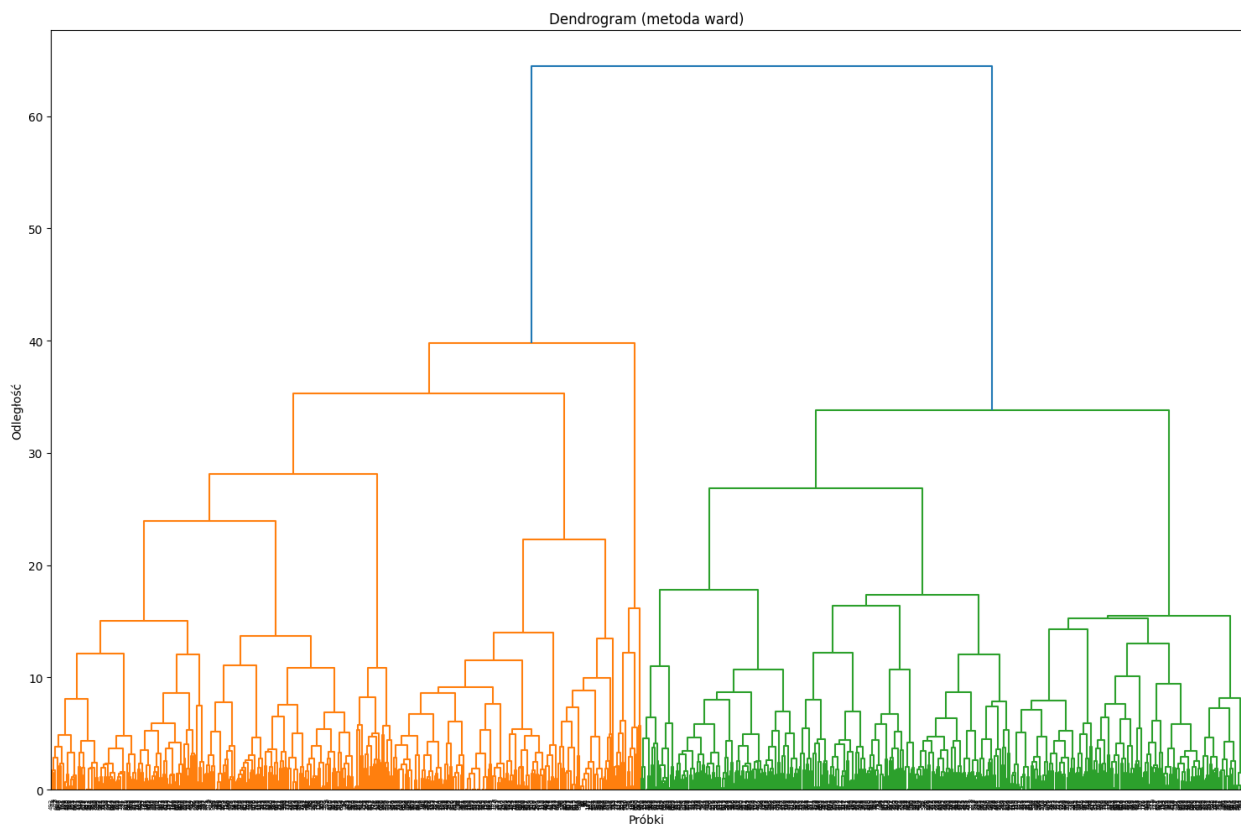
pca = PCA(n_components=2)
dane_zredukowane = pca.fit_transform(dane_standaryzowane)

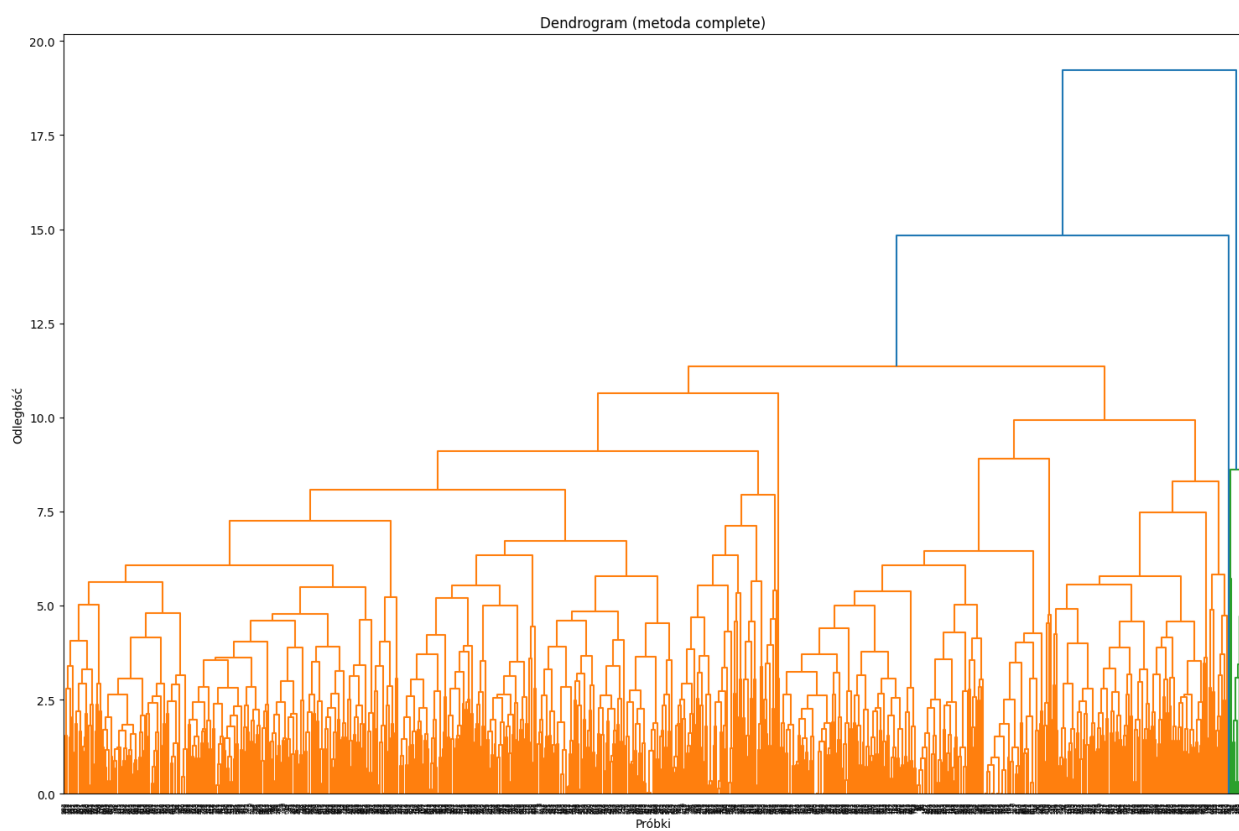
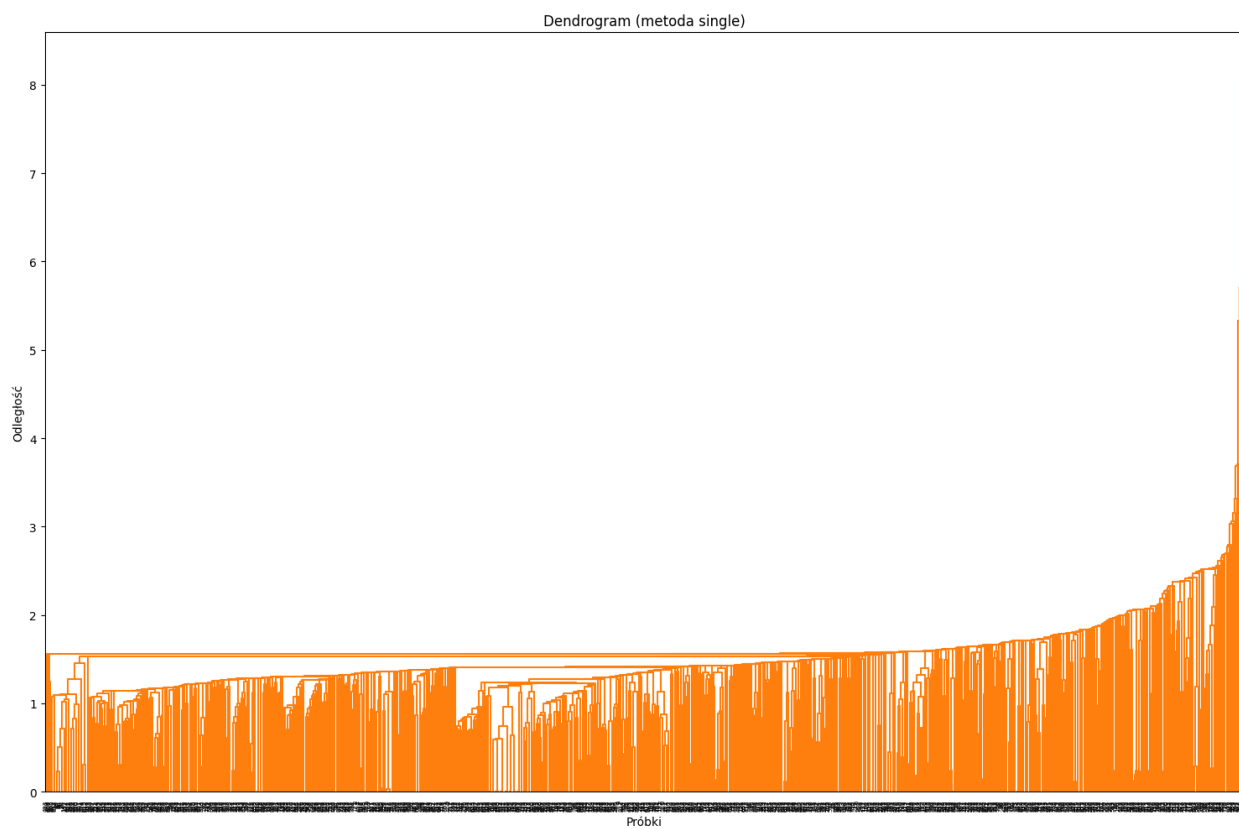
plt.scatter(dane_zredukowane[:, 0], dane_zredukowane[:, 1],
c=klastry_hierarchiczne, cmap='viridis', s=50)
plt.title('Klasteryzacja hierarchiczna na danych zredukowanych PCA')
plt.xlabel('Główna składowa 1')
plt.ylabel('Główna składowa 2')
plt.show()

```

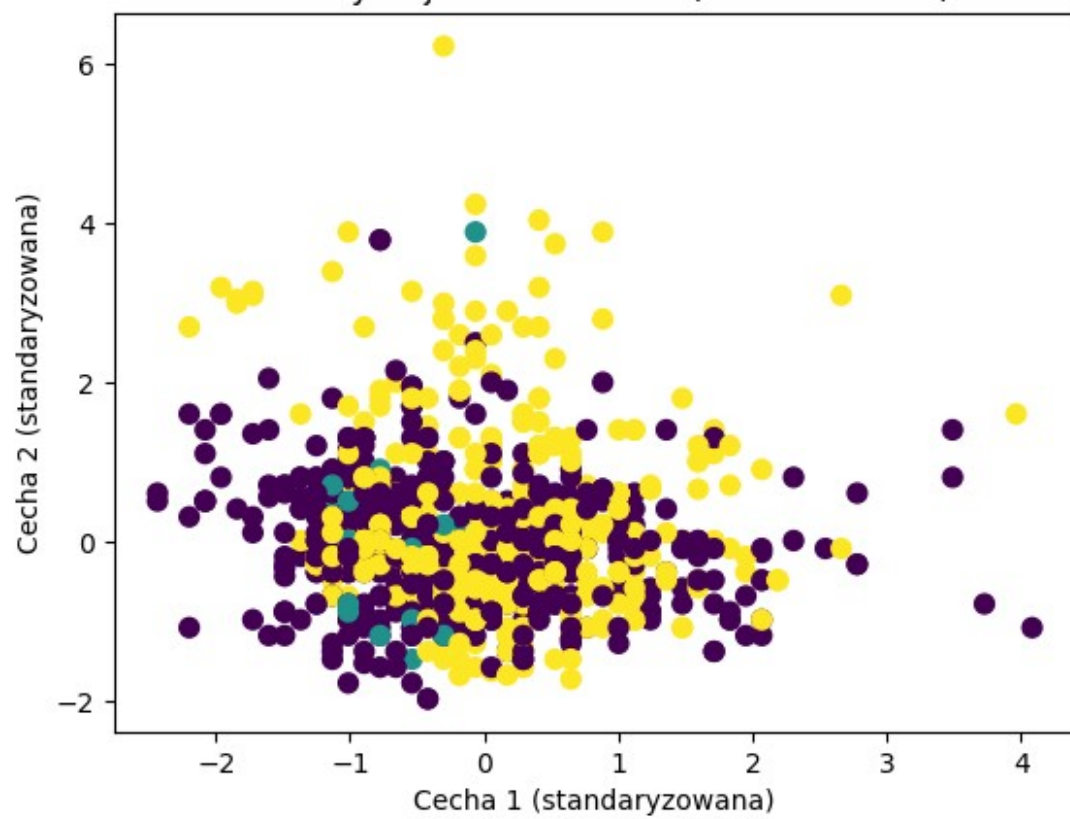


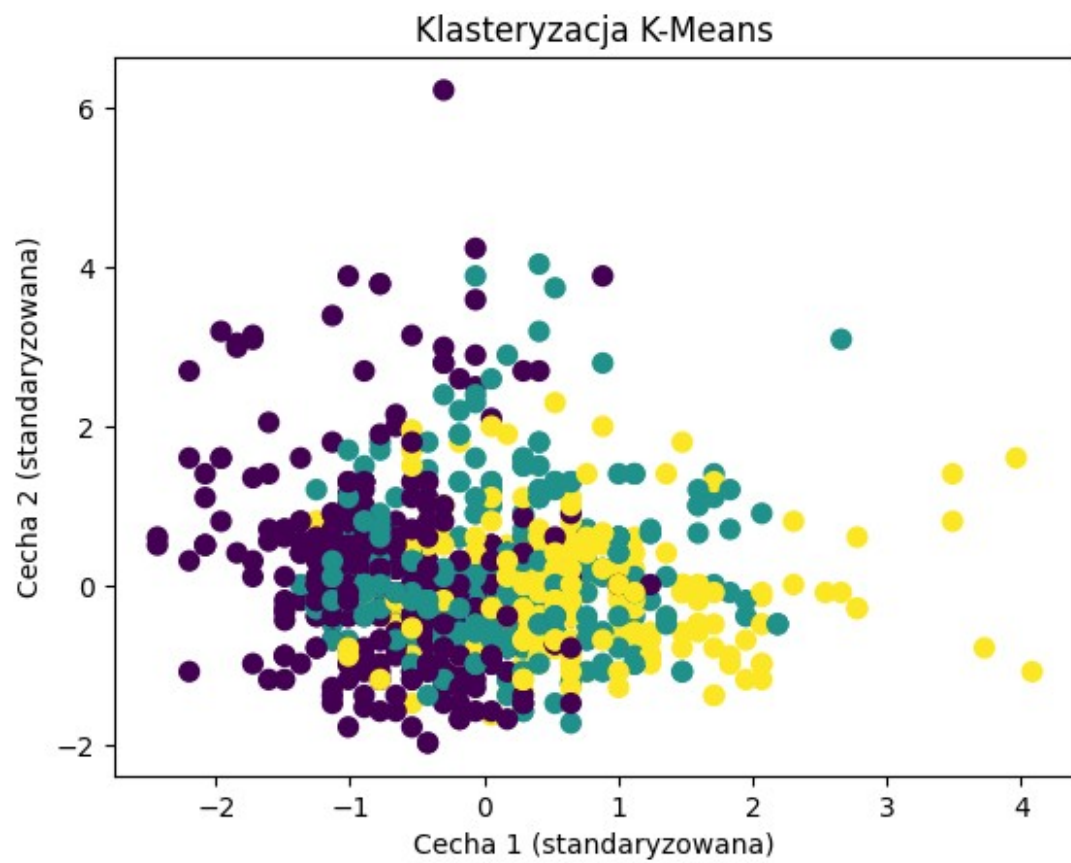
```
plt.scatter(dane_zredukowane[:, 0], dane_zredukowane[:, 1],
c=klastry_kmeans, cmap='viridis', s=50)
plt.title('Klasteryzacja K-Means na danych zredukowanych PCA')
plt.xlabel('Główna składowa 1')
plt.ylabel('Główna składowa 2')
plt.show()
```

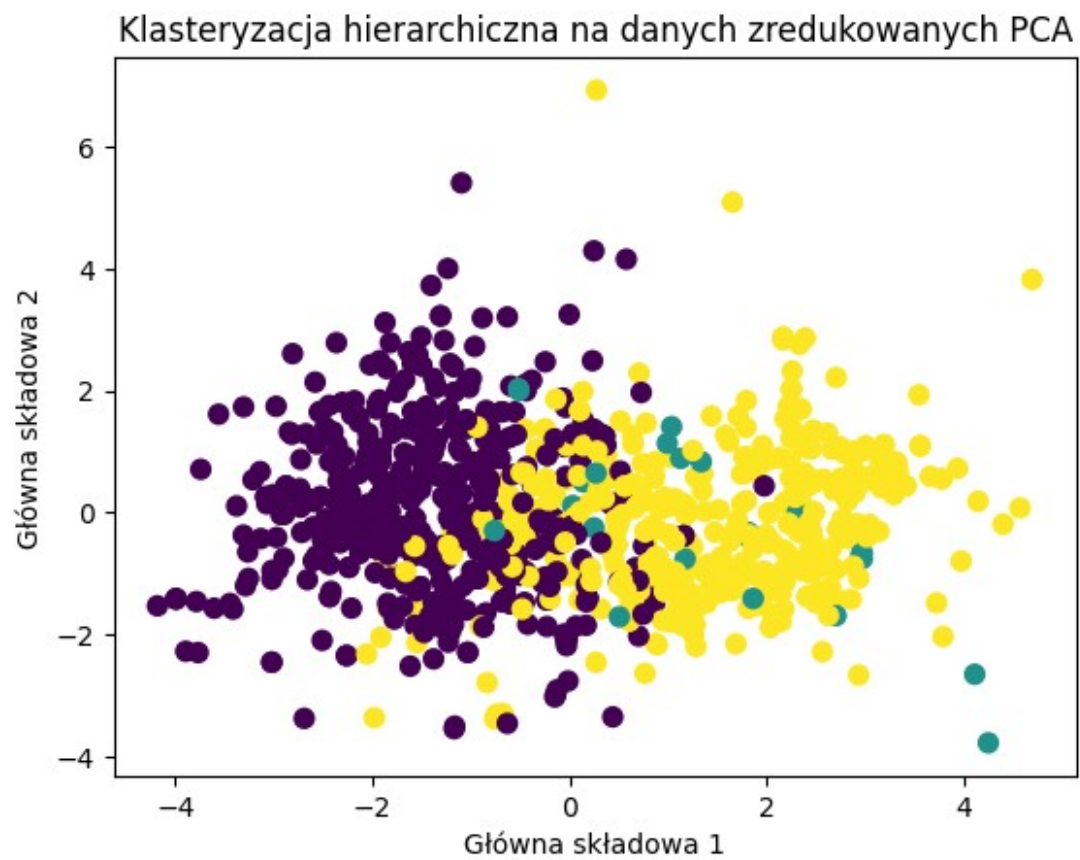




Klasteryzacja hierarchiczna (metoda Warda)







Klasteryzacja K-Means na danych zredukowanych PCA

