# Comparison of Strategies to Inform K-means Clustering

Chris Taylor

DTSA 5510

# Frame the Problem

- K-means requires selection of k before implementing the algorithm
- Two popular methods for selecting k
    - Elbow Plot
        - "Coarse" according to Géron (2021)
    - Silhouette Score
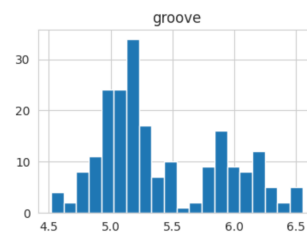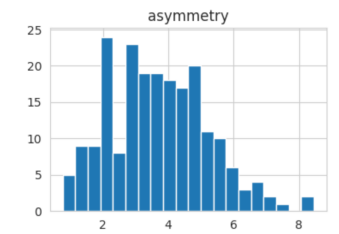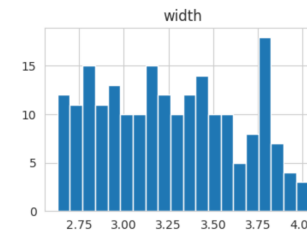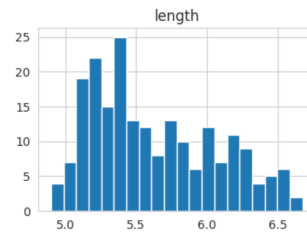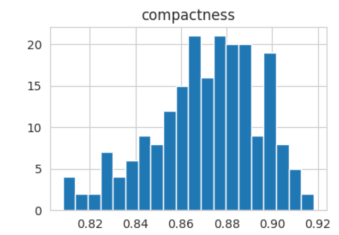- How does the use of each method impact the algorithm?
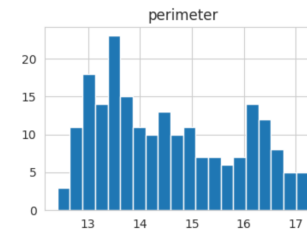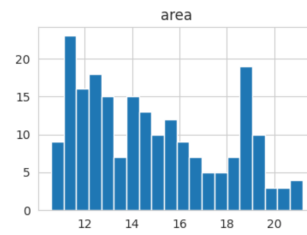
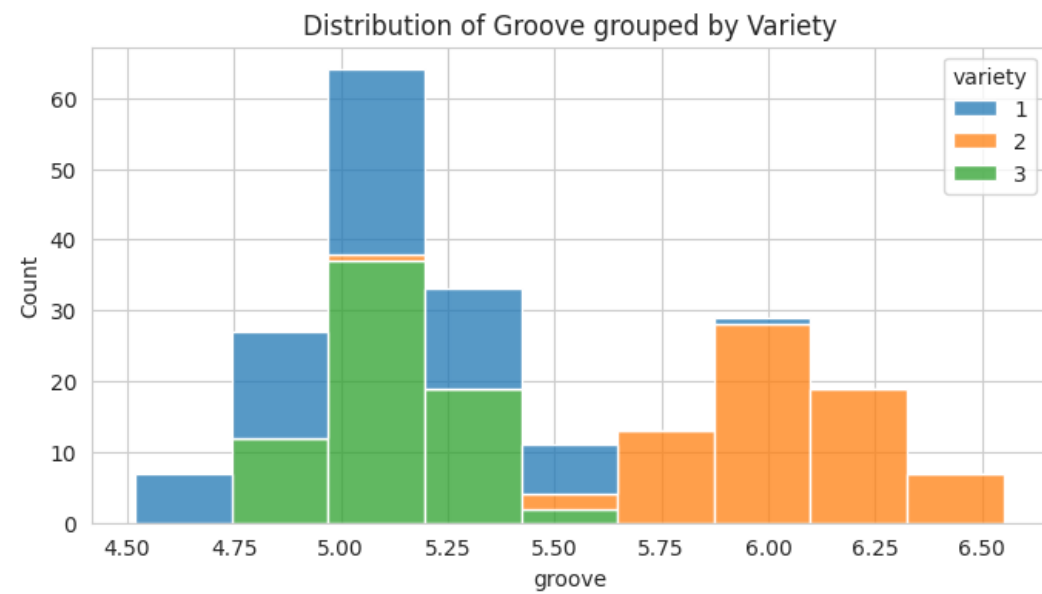# Exploratory Data Analysis

# Description of Data

- 210 observations
- Seven geometric properties
  - Area
  - Perimeter
  - Compactness
  - Kernel length
  - Kernel width
  - Asymmetry coefficient
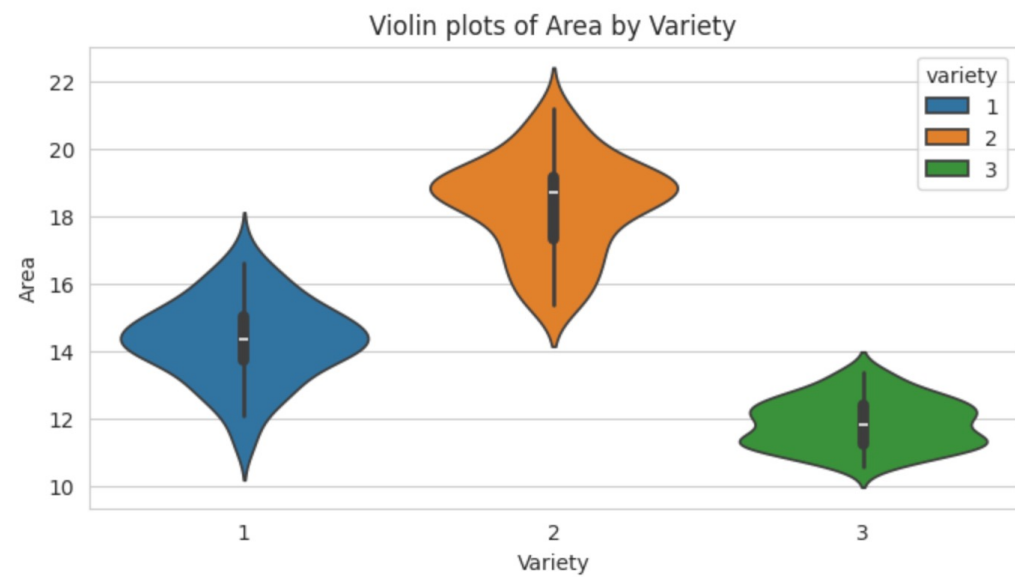  - Kernel groove length

Histogram of Features
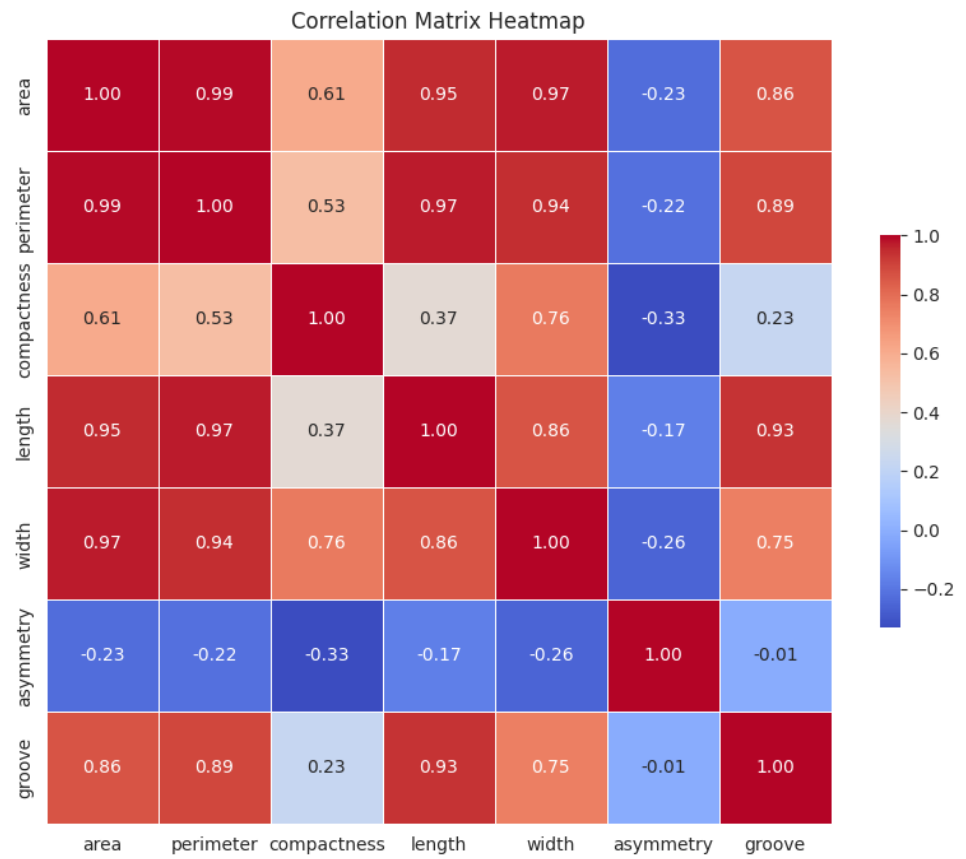
# Histogram of Features



Distribution of Groove grouped by Variety

# Violin Plot of Features



Violin plots of Area by Variety

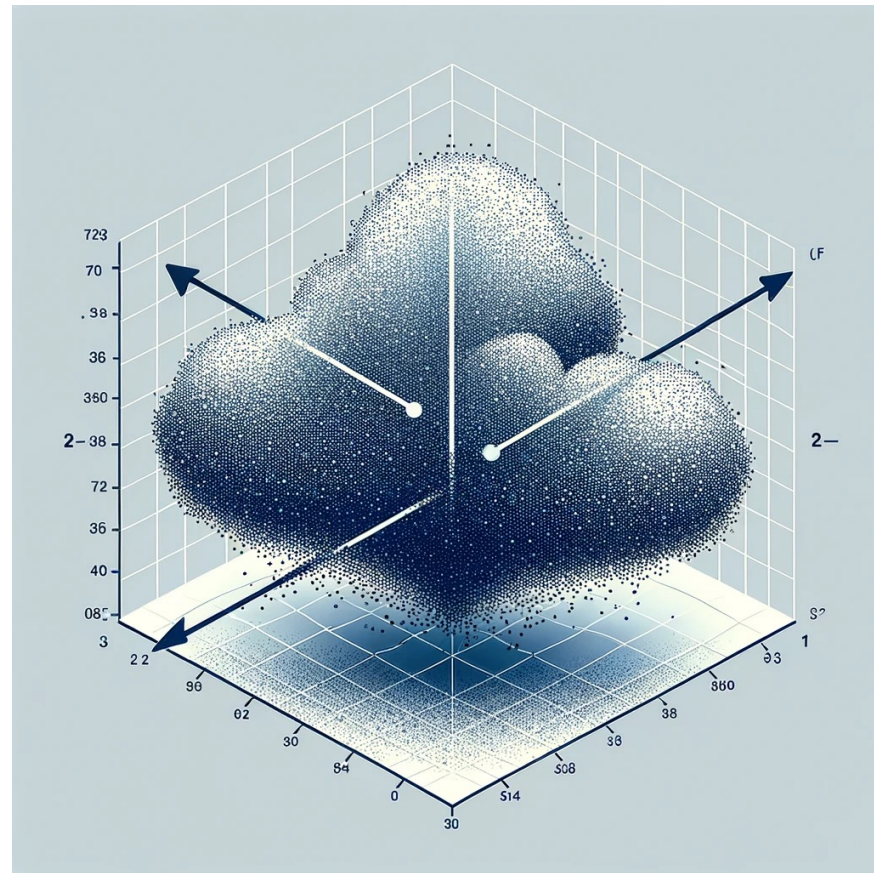Correlation Between Features

Correlation Matrix Heatmap

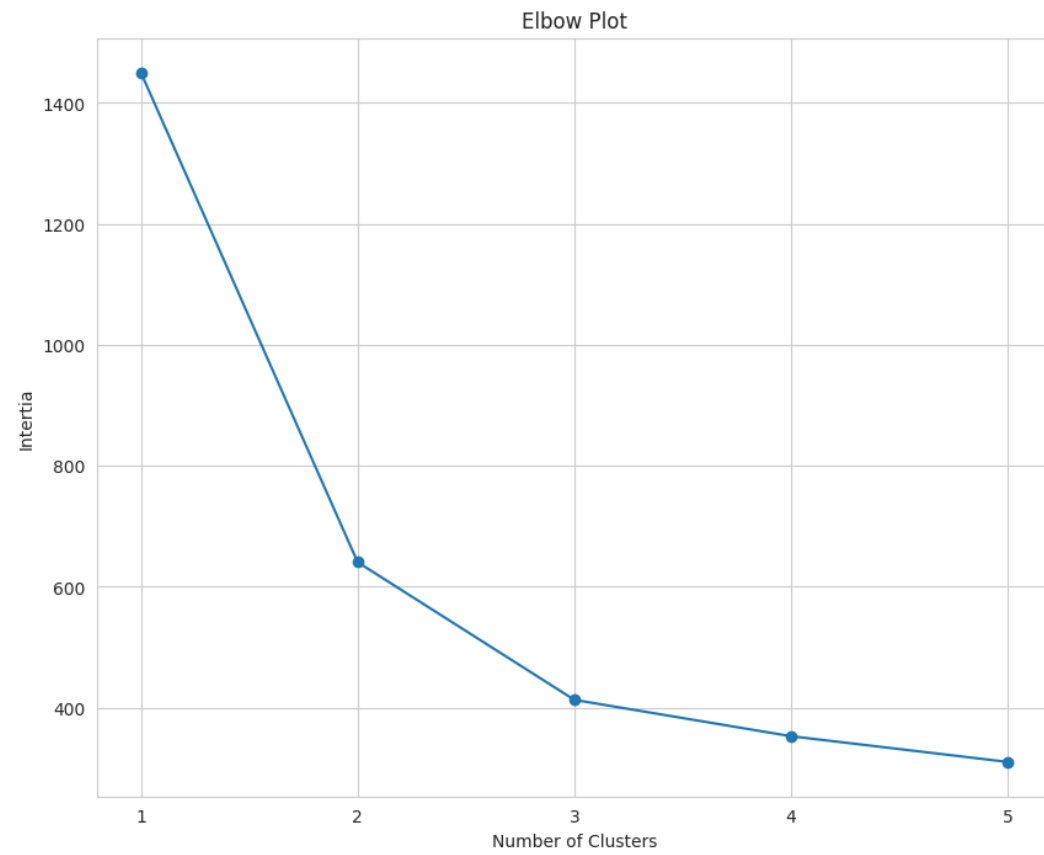# Data Transformation

$$z = \frac{(x - \mu)}{\sigma}$$

# Principle Component Analysis

- Target variance = 0.95
- Reduced dataset used for both implementations of k-means
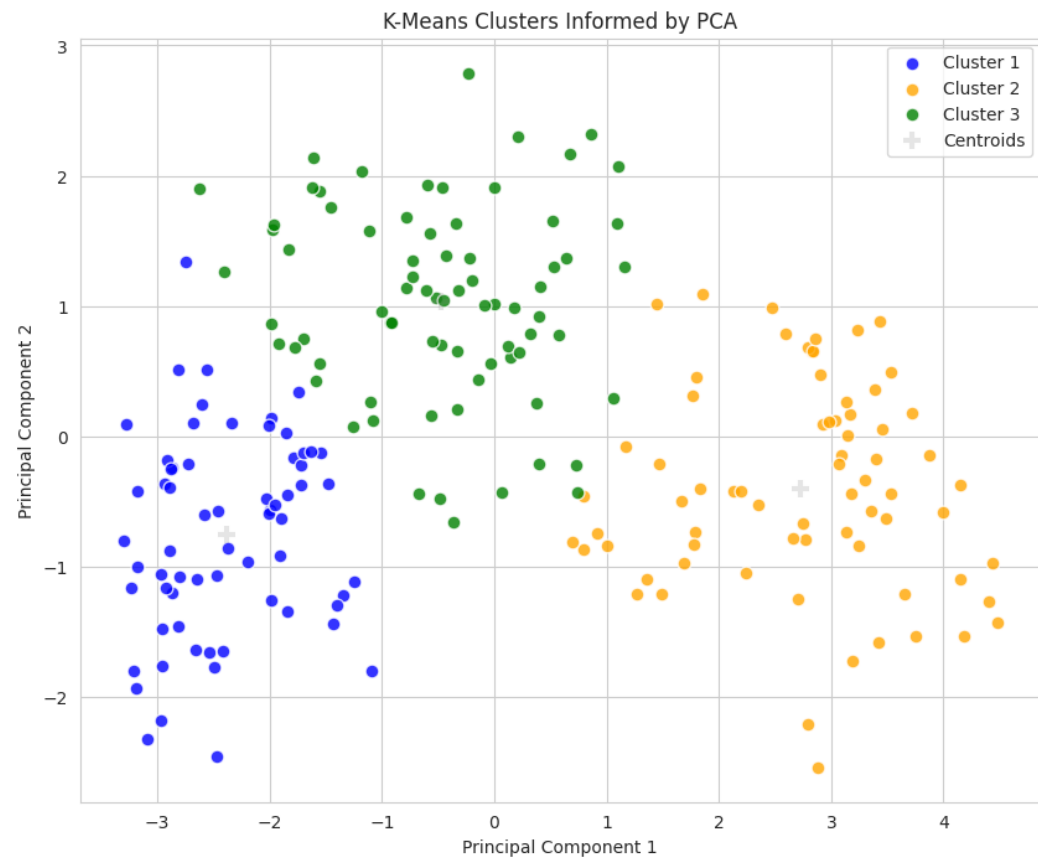
# Elbow Plot Informed K-means Clustering
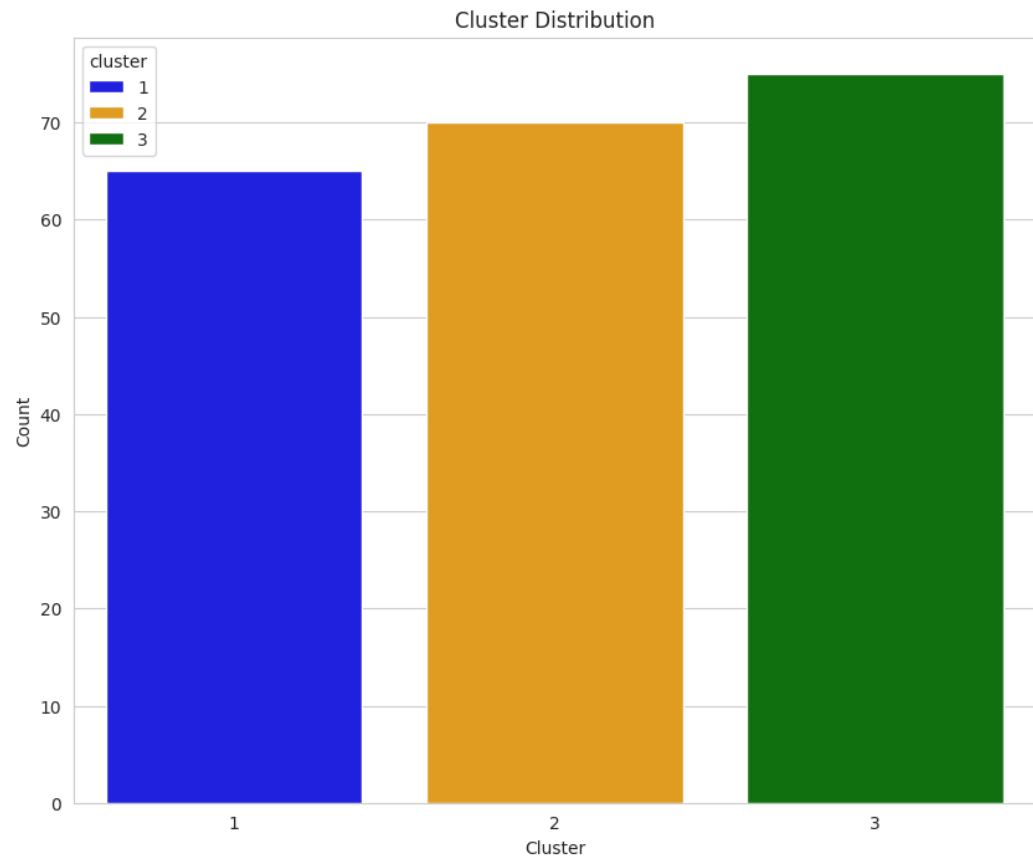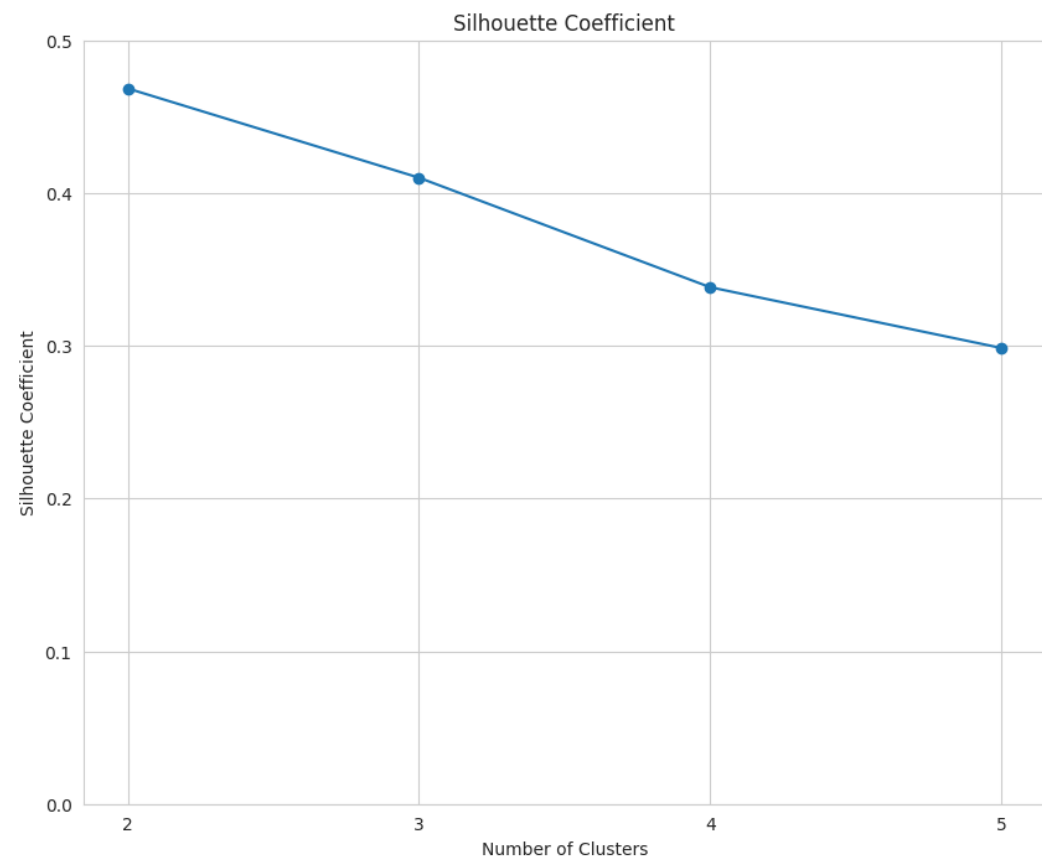
# Elbow Plot

Plot of Clusters

K-Means Clusters Informed by PCA

# Size of Clusters

- Cluster 1 = 65
- Cluster 2 = 70
- Cluster 3 = 75

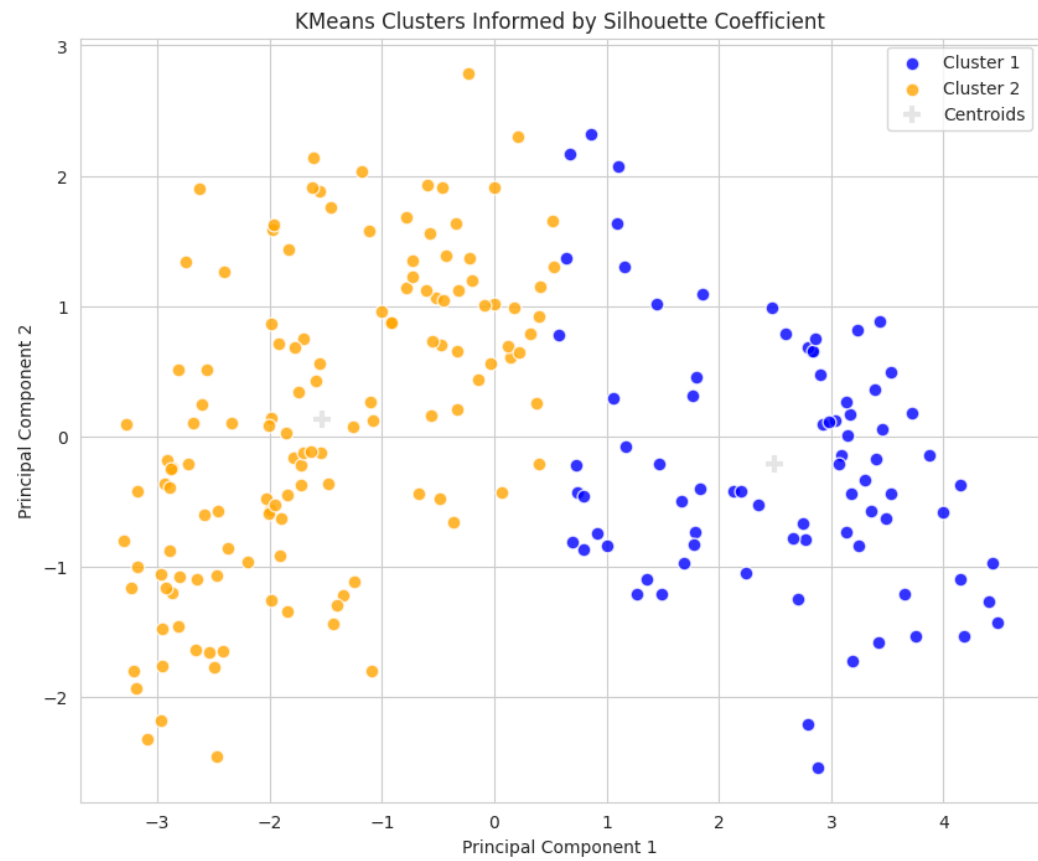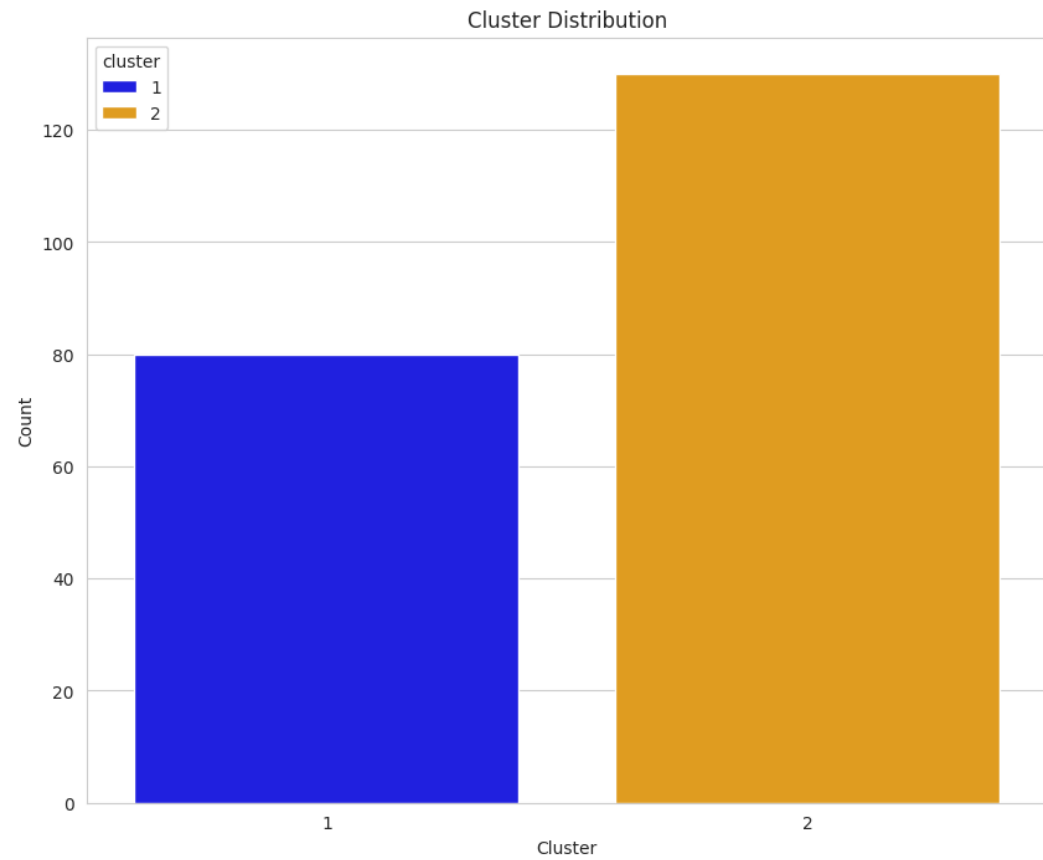# Silhouette Score Informed K-means Clustering

Plot of Clusters

# Size of Clusters

- Cluster 1 = 80
- Cluster 2 = 130



Cluster Distribution

# Comparison of K-means Implementations

# Evaluation Metrics

- Two of three metrics favor K-means informed by Silhouette Score

- Adjusted Rand Index incorporates ground truth

|   | Metric | Elbow Plot | Silhouette Score | Favors |
|---|--------|-----------|------------------|--------|
| 0 | Davies–Bouldin Index | 0.891967 | 0.794722 | SS |
| 1 | Calinski–Harabasz Index | 259.837 | 262.837 | SS |
| 2 | Adjusted Rand Index | 0.773025 | 0.507477 | EP |

# Conclusion

# Conclusion

- K-means informed by Elbow Plot performed better

- Best to explore all options instead of relying on general suggestions

- Future iterations of project could incorporate different dimension reduction techniques