

Introduction to Machine Learning and Clustering – K means - Summary

Key Sections and Concepts

1. Introduction to Machine Learning and Clustering

- **Machine Learning (ML):** The study of algorithms that enable computers to learn from and make decisions based on data.
- **Clustering:** A type of unsupervised learning where the goal is to group similar data points together.

2. Why K-means is Popular

- **Simplicity and Speed:** Despite theoretical limitations in efficiency and quality, K-means remains popular due to its simplicity and ease of implementation, especially for large datasets.
- **Scalability:** K-means can be easily scaled to handle massive datasets through its iterative nature.

3. K-means++ Initialization

- **Improved Initialization:** K-means++ improves the standard K-means algorithm by selecting the first center randomly and each subsequent center with a probability proportional to its contribution to the overall error. This method enhances the chances of finding a global optimum.

4. Determining the Optimal Number of Clusters (k)

- **Silhouette Method:** Measures how similar an object is to its own cluster compared to other clusters. The silhouette score is highest at the optimal k.
- **Elbow Method:** Plots the Within-Cluster-Sum of Squared Errors (WSS) for different k values and identifies the optimal k at the "elbow" point where the WSS starts to diminish.

5. Feature Extraction

- **Importance:** Selecting discriminating and independent features is crucial for any ML algorithm. Often, only a small percentage of measured features carry useful information.
- **Curse of Dimensionality:** High-dimensional data can degrade model performance, necessitating dimensionality reduction techniques.

6. Principal Component Analysis (PCA)

- **Purpose:** PCA is a technique used to reduce the dimensionality of datasets by transforming them into a set of linearly uncorrelated variables called principal components.
- **Process:**
 - Identifies directions of maximum variance.

- Projects data onto a smaller subspace while retaining most of the information.
- Built on eigenvector and eigenvalues concepts, creating a projection matrix to transform the dataset.

Practical Applications

- **Ghelgheli's Teahouse:**
 - **Scenario:** Using clustering algorithms to analyze customer data and optimize business decisions.
 - **Steps:**
 - Data collection and cleaning.
 - Exploratory data analysis (EDA).
 - Application of K-means for clustering customer data based on purchase history, demographics, and feedback.
 - Analysis of clusters to inform strategic decisions.

Conclusion

- **Summary:**
 - K-means remains a popular clustering algorithm due to its simplicity and efficiency.
 - Proper initialization (K-means++) and methods to determine the optimal number of clusters (Silhouette and Elbow methods) can significantly improve clustering results.
 - Feature extraction and dimensionality reduction, particularly PCA, are crucial in handling high-dimensional data effectively.

Introduction to Machine Learning and Clustering

Mahdi KHODADADZADEH

Assistant Professor

Faculty of Geo-Information Science and Earth Observation (ITC)

Department of Geo-information Processing (GIP)

m.khodadadzadeh@utwente.nl

May 2024

Why Kmeans is the most popular algorithm?

- From a theoretical standpoint, k-means is not a good clustering algorithm in terms of efficiency or quality (the running time, locally optimal solution, initialization, ...)
- Why is it one of the top 10 algorithms in data mining? Why is it still popular even as datasets have grown in size?
 - The advantage of k-means is its simplicity. In practice the speed and simplicity of k-means cannot be beat.
 - Scaling k-means to massive data is relatively easy due to its simple iterative nature.
 - Many works have focused on improving this algorithm.

A better way to initialize

- Choosing the centers **one by one in a controlled fashion.**
- **k-means++** algorithm selects only the first center uniformly at random from the data.
- Each subsequent center is selected with a probability proportional to its contribution to the overall error given the previous selections.

Algorithm 1 k -means++(k) initialization.

```
1:  $\mathcal{C} \leftarrow$  sample a point uniformly at random from  $X$ 
2: while  $|\mathcal{C}| < k$  do
3:   Sample  $x \in X$  with probability  $\frac{d^2(x, \mathcal{C})}{\phi_X(\mathcal{C})}$ 
4:    $\mathcal{C} \leftarrow \mathcal{C} \cup \{x\}$ 
5: end while
```

Optimal k value

You may never find the right number of clusters but you can try to find an optimal one!

Run the algorithm for several consecutive number of clusters ($k=1,2,3,\dots,N$).

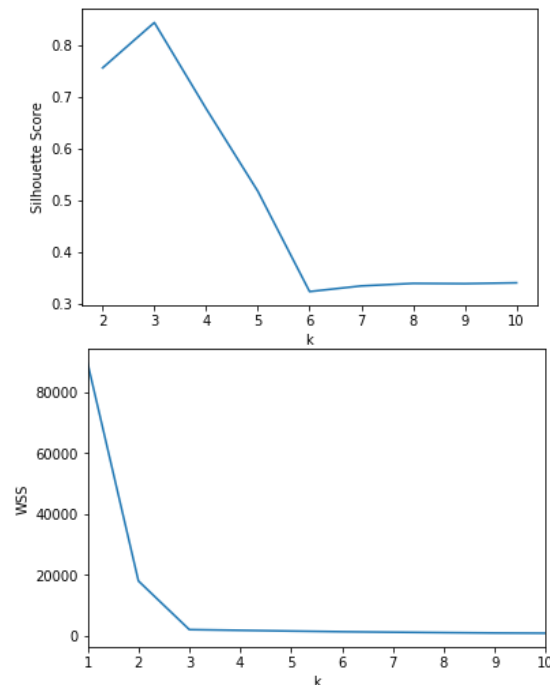
Compute the clustering performance for each number of clusters i.e., k .

Determine the k such that it works well for your problem.

Silhouette Method → measures how similar a point is to its own cluster (cohesion) compared to other clusters (separation) for different values of k . The Silhouette Score reaches its global maximum at the optimal k .

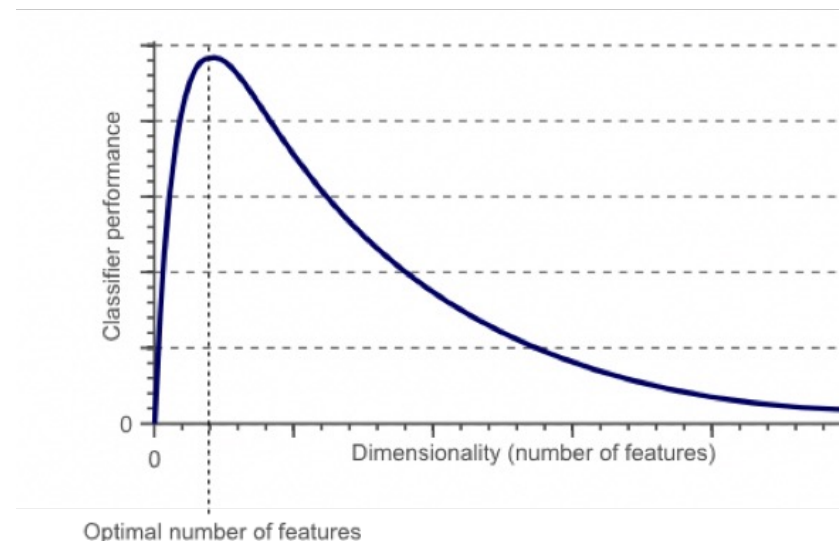
Elbow Method → calculates the Within-Cluster-Sum of Squared Errors (WSS) for different values of k and chooses the k for which WSS becomes first starts to diminish. In the plot of WSS-versus- k , this is visible as an elbow.

Chek out this link: <https://medium.com/analytics-vidhya/how-to-determine-the-optimal-k-for-k-means-708505d204eb>



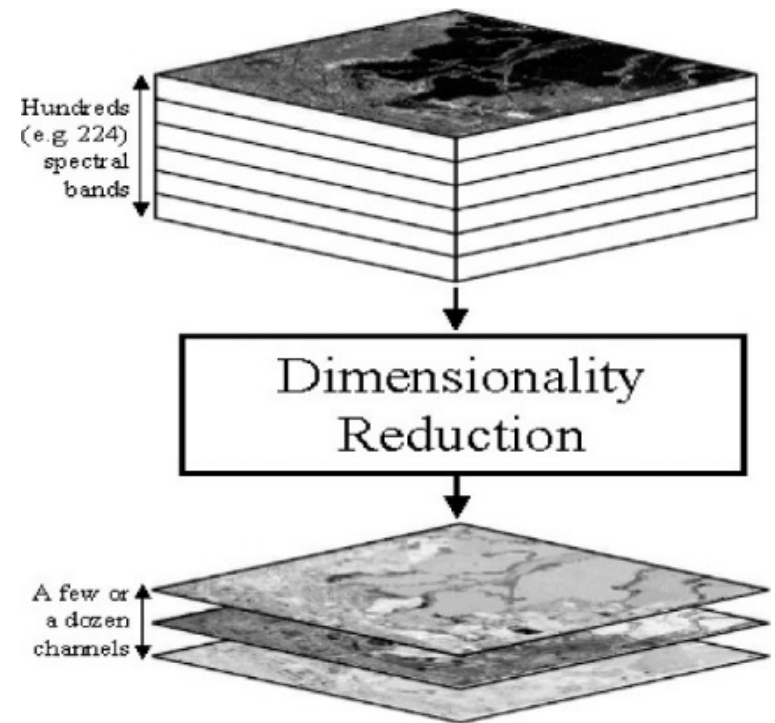
Feature Extraction

- Choosing discriminating and independent features is key to any machine learning algorithm
- In real applications usually many features are measured while only a very small percentage of them carry useful information towards our learning goal
- We usually need an algorithm that compress our feature vector and reduce its dimension
- Curse of dimensionality



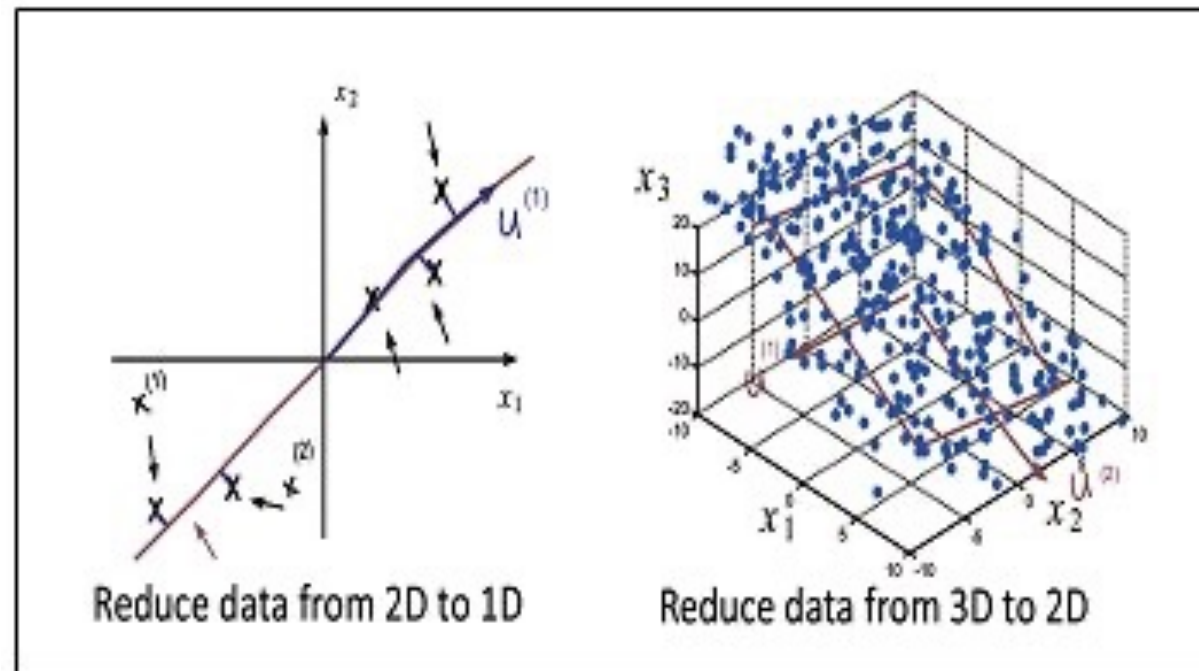
Principal component analysis

- Relatively simple and popular technique
- PCA converts a set of observations into a set of linearly uncorrelated variables, called principal components
- Represents data in a space that better describes the variation
- If a strong correlation between variables exists, the attempt to reduce the dimensionality is reasonable



Principal component analysis

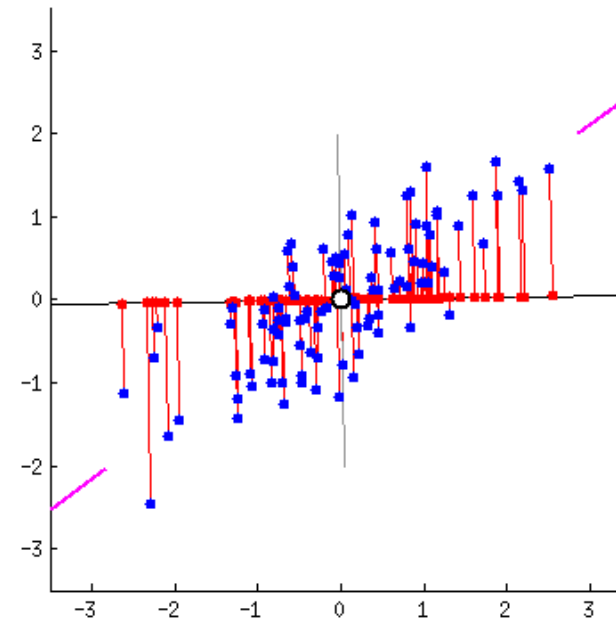
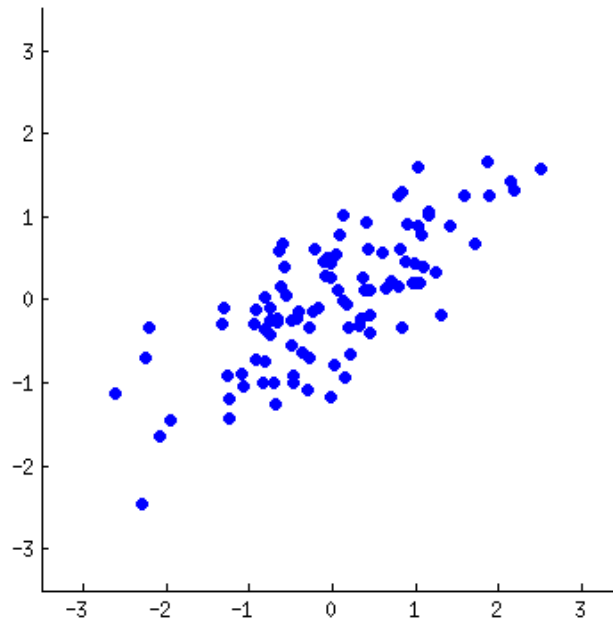
- Identifies directions of maximum variance (in high-dimensional data) and projects the data onto a smaller dimensional subspace while retaining most of the information.
- PCA projects the entire dataset onto a different feature (sub)space



From: Gosavi, V., A. Deshmane, and G. Sable. "Adaptive Neuro Fuzzy Inference System for Facial Recognition." *IOSR Journal of Electrical and Electronics Engineering* 14.3 (2019): 15-22.

Principal component analysis

- What the projections look like for different lines (red dots are projections of the blue dots)
- The reconstruction error are given by the length of the connecting red line



Principal component analysis

- PCA is built on the concepts of eigenvector and eigenvalues
- Creates a projection matrix of the selected k eigenvectors.
- Transforms the original dataset X via the projection matrix and obtains a k -dimensional feature subspace Y

See this link for more details:

<https://towardsdatascience.com/the-mathematics-behind-principal-component-analysis-fff2d7f4b643>