

Unsupervised Kernel Methods

I. Santamaría, S. Van Vaerenbergh

GTAS, Universidad de Cantabria

21 de febrero de 2022

Maître Universitario Oficial **Data Science**



Contents

Introduction

KPCA

Introduction

Kernel clustering

Kernel k-means

Spectral clustering

Conclusions

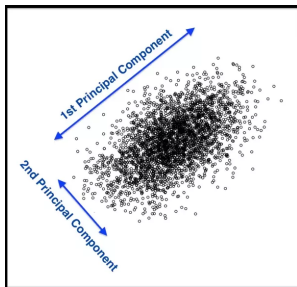
Introduction

Unsupervised learning kernel methods:

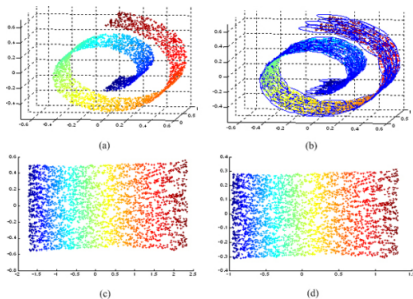
- ▶ Kernel methods for nonlinear dimensionality reduction:
Kernel-PCA (KPCA)
 - ▶ An alternative to other nonlinear dimensionality reduction techniques discussed in M1966 (Data Mining course): LLE, Isomap, t-SNE,...
- ▶ Kernel methods for clustering: **Spectral Clustering/Kernel k-means**

Dimensionality reduction

Linear (PCA)



Nonlinear (Isomap)



PCA (reminder)

- ▶ Normalized input data: $\mathbf{x}_i \in \mathcal{R}^d$ ($i = 1, \dots, n$) (zero-mean unit variance features)

$$\mathbf{X} = [\mathbf{x}_1 \quad \dots, \mathbf{x}_n] \in \mathcal{R}^{d \times n}$$

- ▶ **Sample covariance** matrix ($d \times d$)

$$\mathbf{C} = \frac{1}{n} \mathbf{X} \mathbf{X}^T$$

- ▶ PCA problem (1st component)

$$\max \mathbf{u}_1^T \mathbf{C} \mathbf{u}_1 \quad \text{s.t.} \quad \|\mathbf{u}_1\|_2^2 = 1$$

- ▶ Solution: main eigenvector of \mathbf{C}

$$\mathbf{C} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^T \quad \mathbf{U} = [\mathbf{u}_1 \quad \dots \quad \mathbf{u}_d]$$

- ▶ To obtain the first r principal components, $\mathbf{U}_r = [\mathbf{u}_1, \dots, \mathbf{u}_r]$ ($n \times r$), the problem amounts to

$$\max \operatorname{tr}(\mathbf{U}_r^T \mathbf{C} \mathbf{U}_r), \quad \text{s.t.} \quad \mathbf{U}_r^T \mathbf{U}_r = \mathbf{I},$$

whose solution is

$$\mathbf{C} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^T \quad \mathbf{U}_r = [\mathbf{u}_1 \quad \dots \quad \mathbf{u}_r \quad \mathbf{u}_{r+1} \quad \dots \quad \mathbf{u}_d]$$

- ▶ PCA can also be solved starting from the kernel matrix $\mathbf{K} = \mathbf{X}^T \mathbf{X}$

KPCA

- ▶ KPCA \rightarrow PCA in the transformed (feature) space
- ▶ Let $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ be a dataset in the input space
- ▶ We want to find maximum variance projections of the transformed vectors $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$
- ▶ The first principal component of $\Phi(\mathbf{x})$ can be written as

$$y_1 = \sum_{i=1}^n \alpha_{1,i} k(\mathbf{x}, \mathbf{x}_i) = \mathbf{k}_i^T \boldsymbol{\alpha}_1$$

where $\boldsymbol{\alpha}_1 = \lambda_1^{-1/2} \mathbf{v}_1$ is an $n \times 1$ vector, \mathbf{v}_1 is the largest eigenvector of the kernel matrix, \mathbf{K} , and λ_1 is the corresponding eigenvalue

- ▶ Subsequent principal components are obtained similarly

Centering

- It is common practice to apply KPCA over zero-mean data (in the feature space)

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{0} \nRightarrow \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i) = \boldsymbol{\mu} = \mathbf{0}$$

- Centering or mean removal in the feature space

$$\Phi_c(\mathbf{x}) = \Phi(\mathbf{x}) - \boldsymbol{\mu}$$

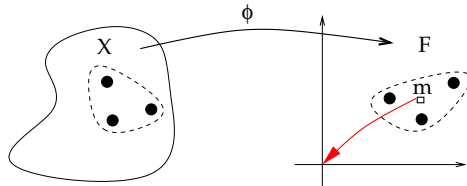
- The centered kernel matrix is

$$\begin{aligned} k_c(\mathbf{x}, \mathbf{y}) &= k(\mathbf{x}, \mathbf{y}) - \frac{1}{n} \sum_i k(\mathbf{x}, \mathbf{x}_i) - \frac{1}{n} \sum_i k(\mathbf{y}, \mathbf{x}_i) \\ &\quad + \frac{1}{n^2} \sum_i \sum_j k(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

- The centered **kernel matrix** can be written as

$$\mathbf{K}_c = \mathbf{K} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{K} - \frac{1}{n} \mathbf{K} \mathbf{1} \mathbf{1}^T + \frac{1}{n^2} \mathbf{1} \mathbf{1}^T \mathbf{K} \mathbf{1} \mathbf{1}^T = \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right) \mathbf{K} \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^T \right)$$

where $\mathbf{1} = [1, \dots, 1]^T$ and \mathbf{I} is the identity matrix



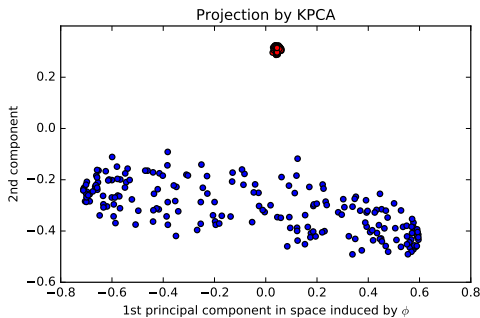
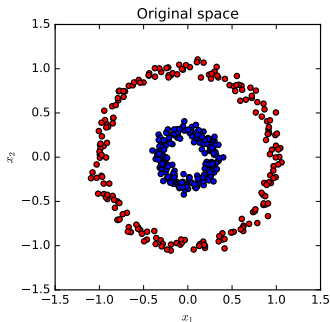
- ▶ Input: Data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{R}^d$, number of principal components or projections, r , kernel parameters (σ^2 or γ)
- ▶ Output: $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathcal{R}^r$

KPCA

1. Compute the kernel matrix \mathbf{K}
2. Kernel matrix centering: $\mathbf{K} = (\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^T) \mathbf{K} (\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^T)$
3. $[\mathbf{V}, \mathbf{\Lambda}] = \text{eig}(\mathbf{K})$
4. $\alpha_j = \lambda_j^{-1/2} \mathbf{v}_j, j = 1, \dots, r$
5. for $i = 1 : n$
 - ▶ $\mathbf{k}_i = [k(\mathbf{x}_i, \mathbf{x}_1) \quad \dots \quad k(\mathbf{x}_i, \mathbf{x}_n)]^T$
 - ▶ $\mathbf{y}_i = [\alpha_1^T \mathbf{k}_i \quad \dots \quad \alpha_r^T \mathbf{k}_i]^T$

Example

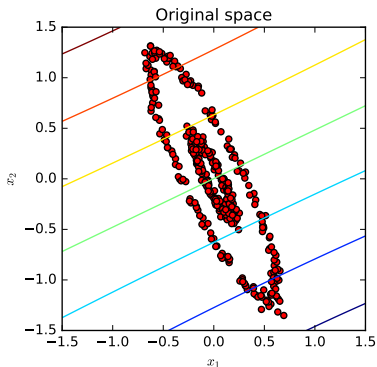
KPCA can make data linearly separable



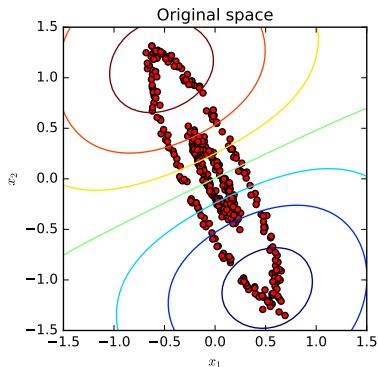
Example

KPCA can extract nonlinear correlations in the dataset

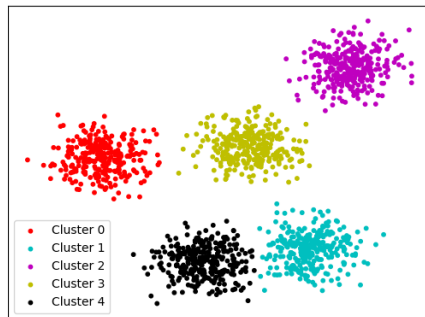
$\gamma = 0.01$



$\gamma = 0.5$



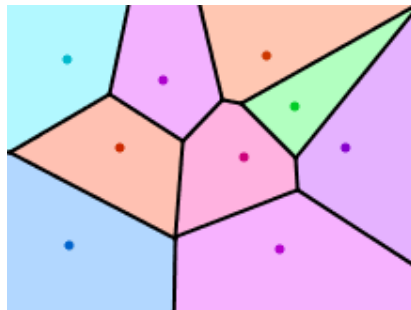
Kernel clustering



- ▶ Kernel methods can also be applied to **clustering**
- ▶ Two popular kernel-based clustering methods are
 1. Kernel k-means
 2. Spectral clustering

k-means in the input space

- ▶ **k-means** is probably the most popular clustering method
- ▶ Input: Data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathcal{R}^p$, and number of clusters k
- ▶ Output: k centroids $\mu_1, \dots, \mu_k \in \mathcal{R}^p$
- ▶ The centroids split the input space into k disjoint Voronoi regions or clusters



- ▶ The clustering problem is to find the optimal centroids that minimize a distortion criterion

$$D(\mu_1, \dots, \mu_k) = \sum_{j=1}^k \sum_{\mathbf{x}_n \in \mathcal{C}_j} \|\mathbf{x}_n - \mu_j\|_2^2$$

- ▶ Each cluster, \mathcal{C}_j , is defined by its corresponding centroid μ_j
- ▶ To solve the problem we have to:
 - ▶ Assign patterns to clusters $\mathbf{x}_n \rightarrow \mathcal{C}_j$
 - ▶ Estimate centroids μ_j
- ▶ There is no closed-form solution, so we have to resort to iterative algorithms

k-means

1. Random initialization of centroids $\mu_j \in \mathcal{R}^p, j = 1, \dots, k$
2. **Assign patterns to clusters/centroids**: assign each pattern \mathbf{x}_n to its closest centroid

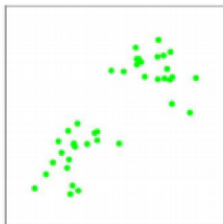
$$\mathbf{x}_n \in \mathcal{C}_i, \quad i = \underset{j=1, \dots, k}{\operatorname{argmin}} \|\mathbf{x}_n - \mu_j\|_2^2$$

3. **Update centroids** as

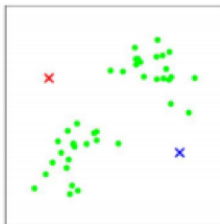
$$\mu_i = \frac{1}{n_i} \sum_{\mathbf{x}_n \in \mathcal{C}_i} \mathbf{x}_n$$

Monotonic convergence, possibly to a local minimum!

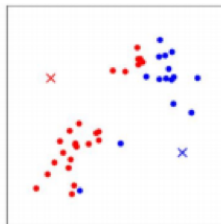
Example



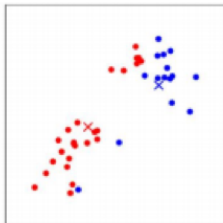
(a)



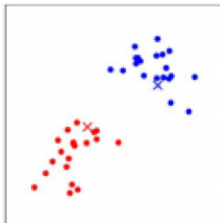
(b)



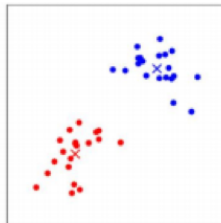
(c)



(d)



(e)



(f)

Kernel k-means

- ▶ The “kernelized” version of the algorithm applies k-means in the feature space
- ▶ Clustering problem: to find centroids/clusters that minimize

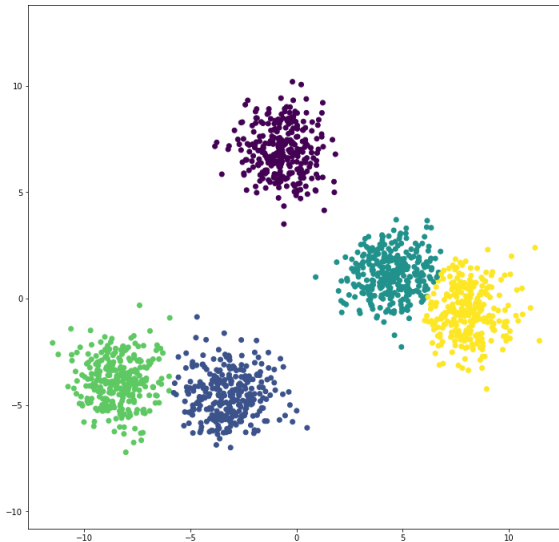
$$D(\mu_1, \dots, \mu_k) = \sum_{j=1}^k \sum_{\mathbf{x}_n \in \mathcal{C}_j} \|\Phi(\mathbf{x}_n) - \mu_j\|_2^2$$

- Distances can be written in terms of the kernel function as

$$\begin{aligned}\|\Phi(\mathbf{x}_n) - \mu_j\|_2^2 &= \|\Phi(\mathbf{x}_n) - \frac{1}{n_j} \sum_{j \in \mathcal{C}_j} \Phi(\mathbf{x}_j)\|_2^2 \\ &= k(\mathbf{x}_n, \mathbf{x}_n) - \frac{2}{n_j} \sum_{j \in \mathcal{C}_j} k(\mathbf{x}_n, \mathbf{x}_j) + \frac{1}{n_j n_i} \sum_{j \in \mathcal{C}_j} \sum_{i \in \mathcal{C}_i} k(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

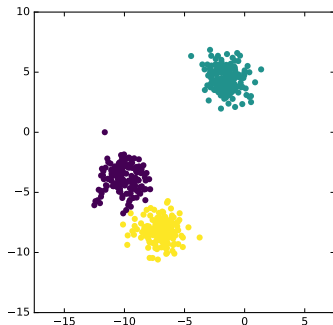
therefore, the k-means algorithm can directly be applied in the feature space

Example

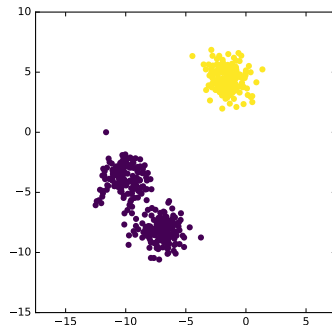


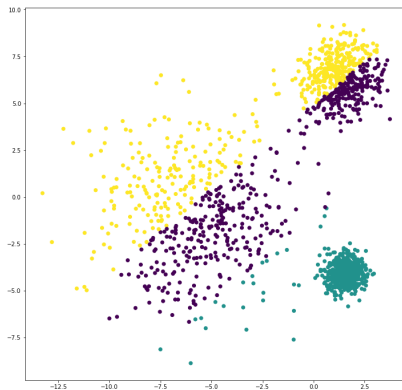
k-means assumes the number of clusters to be known

$k = 3$



$k = 2$

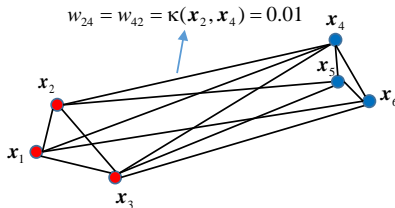


[illegible]

Spectral clustering

- ▶ A popular kernel method for clustering
- ▶ One intuitive way to understand spectral clustering is as a partition, or cut, of a similarity graph defined by the kernel matrix
- ▶ Given a set of patterns $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, we define an undirected graph as $G = (V, E)$ where
 - ▶ The patterns are the nodes or graph vertices V ,
 - ▶ All nodes are connected through edges (fully connected graph)
 - ▶ The weight of an edge between two patterns measures the similarity between them as: $w_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$

Similarity graph



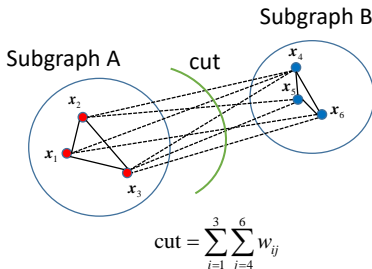
$$W = K = \begin{bmatrix} 1 & 0.9 & 0.8 & 0.05 & 0.01 & 0 \\ 0.9 & 1 & 0.7 & 0.01 & 0.01 & 0.01 \\ 0.8 & 0.7 & 1 & 0.02 & 0.01 & 0.01 \\ 0.05 & 0.01 & 0.02 & 1 & 0.8 & 0.7 \\ 0.01 & 0.01 & 0.01 & 0.8 & 1 & 0.8 \\ 0 & 0.01 & 0.01 & 0.7 & 0.8 & 1 \end{bmatrix}$$

If the clusters are well separated, \mathbf{K} is (approximately) a block-diagonal matrix

Graph cut

- Assuming there are two clusters, the problem would be to cut the graph into two disjoint subgraphs such that the sum of the edges separating the two subgraphs is minimum

$$\text{cut}(A, B) = \sum_{i \in A, j \in B} w_{ij}$$



- ▶ Spectral clustering solves the graph cut problem
- ▶ It applies PCA to the Laplacian of the graph

$$\mathbf{L} = \mathbf{D} - \mathbf{W}$$

where \mathbf{D} is a diagonal matrix with elements $d_i = \sum_{j=1}^n w_{ij}$

- ▶ The eigenvectors corresponding to the largest k eigenvalues of \mathbf{L} contain information about the k connected subgraphs (clusters)

Spectral Clustering

1. Input:

- ▶ Patterns $(\mathbf{x}_1, \dots, \mathbf{x}_n)$, $\mathbf{x}_i \in \mathcal{R}^d$
- ▶ number of clusters, k
- ▶ Kernel matrix \mathbf{K} with $k(i, j) = \exp(-\gamma|\mathbf{x}_i - \mathbf{x}_j|^2)$

2. Compute the graph Laplacian

$$\mathbf{L} = \mathbf{D} - \mathbf{K}$$

3. Store the first k eigenvectors of \mathbf{L} into matrix

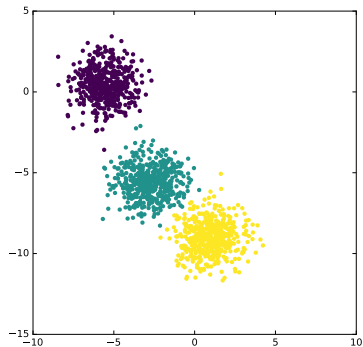
$$\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_k] \in \mathcal{R}^{n \times k}$$

4. Let $\mathbf{y}_i \in \mathcal{R}^k$ ($i = 1, \dots, n$) be the i -th row of \mathbf{V}

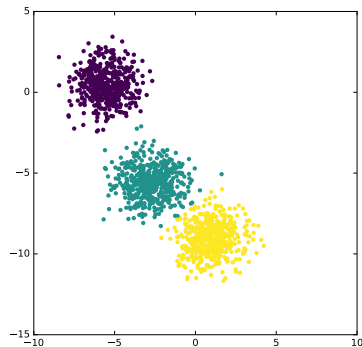
5. Apply k-means to the set of row vectors \mathbf{y}_i , $i = 1, \dots, n$

6. Output: Clusters C_1, \dots, C_k obtained from k-means

Kernel k-means



Spectral Clustering



Conclusions

- ▶ **KPCA**: PCA applied in the feature space
 - ▶ Nonlinear dimensionality reduction
 - ▶ Nonlinear correlation analysis
- ▶ Kernel methods for clustering
 - ▶ **Kernel k-means**: k-mean applied in the feature space
 - ▶ **Spectral clustering**: KPCA + k-means