

# Kernel K-Means

Saturday, March 11, 2017

10:28 AM

(1)  
T

Mercer's Theorem. Let  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ . Let  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  be an inner product on the Hilbert space  $\mathcal{H}$ . There exists a map  $\varphi: \mathcal{X} \rightarrow \mathcal{H}$  such that

$$k(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$$

if and only if for any square integrable function  $g(x)$ , i.e.  $\int |g(x)|^2 dx < \infty$ , the following holds:

$$\iint dx dy g(x) k(x, y) g(y) \geq 0.$$

\_\_\_\_\_ h \_\_\_\_\_

$$\begin{aligned} \|\varphi(x) - \varphi(y)\|_{\mathcal{H}}^2 &= (\varphi(x) - \varphi(y))^T (\varphi(x) - \varphi(y)) \\ &= \varphi(x)^T \varphi(x) - 2 \varphi(x)^T \varphi(y) + \varphi(y)^T \varphi(y) \\ &= k(x, x) - 2 k(x, y) + k(y, y) \end{aligned}$$

↑    //

$D_{\mathcal{H}}(x, y)$

K-means  $E = \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^n z_{ki} \|x_i - m_k\|^2$

where  $z_{ki} = \begin{cases} 1 & \text{if } x_i \in C_k \\ 0 & \text{otherwise} \end{cases}$ ,  $m_k = \frac{\sum_{i=1}^n z_{ki} x_i}{\sum_{i=1}^n z_{ki}} = \frac{1}{|C_k|} \sum_{x_i \in C_k} x_i$

]

②

Kernel k-means do the same thing  
but on  $\mathcal{H}$ .

$$E = \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^m z_{ki} \| \varphi(x_i) - m_k \|_{\mathcal{H}}^2$$

$$m_k = \frac{1}{|C_k|} \sum_{x_i \in C_k} \varphi(x_i) = \frac{\sum_{i=1}^m z_{ki} \varphi(x_i)}{\sum_{i=1}^m z_{ki}}.$$

$$\begin{aligned} \| \varphi(x_i) - m_k \|_{\mathcal{H}}^2 &= \| \varphi(x_i) - \frac{1}{m_k} \sum_{l=1}^m z_{kl} \varphi(x_l) \|_{\mathcal{H}}^2 \\ &= \varphi(x_i)^T \varphi(x_i) - 2 \varphi(x_i)^T \frac{1}{m_k} \sum_{l=1}^m z_{kl} \varphi(x_l) \\ &\quad + \frac{1}{m_k^2} \sum_{l=1}^m \sum_{p=1}^m z_{kl} z_{lp} \varphi(x_l)^T \varphi(x_p) \\ &= k(x_i, x_i) - \frac{2}{m_k} \sum_{l=1}^m k(x_i, x_l) \\ &\quad + \frac{1}{m_k^2} \sum_{l=1}^m \sum_{p=1}^m z_{kl} z_{lp} k(x_l, x_p) \end{aligned}$$

Thus the kernel k-means objective function is

$$E = \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^m z_{ki} \left( k(x_i, x_i) - \frac{2}{m_k} \sum_{l=1}^m z_{kl} k(x_i, x_l) + \frac{1}{m_k^2} \sum_{l=1}^m \sum_{p=1}^m z_{kl} z_{lp} k(x_l, x_p) \right)$$

where  $m_k = \sum_{i=1}^m z_{ki}$ . Only depends on the

Gram matrix  $K_{ij} = k(x_i, x_j)$

We will see that this can be  
written as

$$\text{Tr}(G K G^T) \quad \text{kernel k-means.}$$

Spectral clustering is a weighted kernel  
k-means

$$\text{Tr}(G D^{1/2} K D^{1/2} G^T)$$

Problems:

1. k-means:  $O(n \times d)$   
kernel k-means:  $O(n^2)$
2. In a sufficiently high dimensioned space, every data can be separable. However, a choice of kernel is data/applications dependent, and usually introduces new parameters.
3. The wrong choice of kernel can lead to poor results, even worse than k-means.

Options for scalability issues:  
→ distributed,  
→ does not improve complexity

1. Reordered Kernel k-means
2. Sampling based approximation  
(kernel matrix approx, random projections)

Approximating the kernel:

(4)

Sample  $m \ll n$  points  $\rightarrow \{y_1, \dots, y_m\}$

Compute  $\{k_A \in \mathbb{R}^{m \times m}, (k_A)_{ij} = \varphi(y_i)^T \varphi(y_j)\}$   
 $\{k_B \in \mathbb{R}^{m \times m}, (k_B)_{ij} = \varphi(x_i)^T \varphi(y_j)\}$

$$m_h \approx \sum_{j=1}^m \alpha_{kj} \varphi(y_j)$$

$$\min_{\alpha} \max_k \frac{1}{2} \sum_{h=1}^k \sum_{i=1}^m z_{ih} \| \varphi(x_i) - \sum_{j=1}^m \alpha_{kj} \varphi(y_j) \|^2$$

- This is equivalent to run k-means on  $K_B K_A^{-1} K_B^T$  (Nyström approx.)
- $O(n \cdot m)$  almost linear on  $n$ .
- $O(\frac{1}{m})$  approx. error.

Methods:

- Nyström. Randomly select columns  
 $k = K_B K_A^{-1} K_B^T$
- CUR. Select best columns C and rows R. Find U s.t  
 $\min \|k - CUR\|_F$

(57)

- Find  $N$  such that

$$\hat{h} = h + N$$

is sparse, and  $E[N_{ij}] = 0$  and  $\text{Var}[N_{ij}]$  is small

Kernel PCA Project data into the first  $C$  eigenvectors of  $K$ , cluster on the eigen space.

## Nonlinear Random Projections

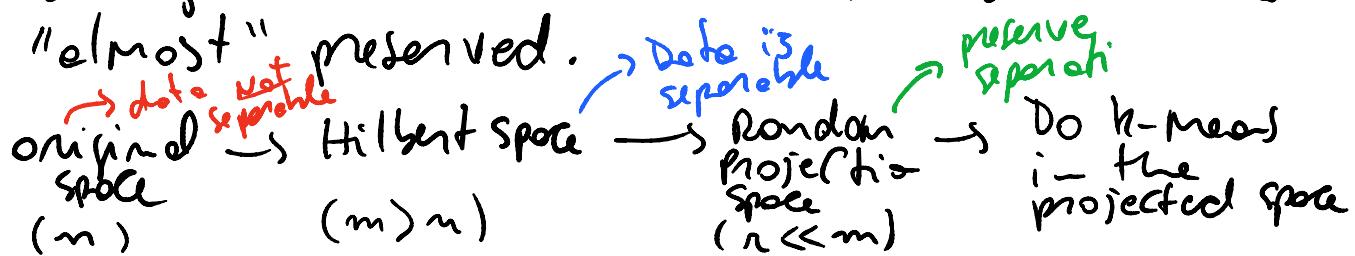
Pick a random matrix  $R \in \mathbb{R}^{d \times r}$ . Transform data from the Hilbert space:

$$x' = \frac{1}{\sqrt{n}} R^T x, \quad y' = \frac{1}{\sqrt{n}} R^T y$$

## Johnson-Lindenstrauss theorem

$$(1 - \epsilon) \|x - y\|^2 \leq \|x' - y'\|^2 \leq (1 + \epsilon) \|x - y\|^2$$

Provided  $R$  is orthogonal or  $r$  is large enough. This means that distances are "almost" preserved.



1

(6)

## Unsupervised Nonparametric Kernel Learning Algorithm. Liu et al (2013)

$$\min_{K, U} \text{Tr}(KL) + \frac{1}{2} U \left( K + \frac{I}{c} KL + \frac{I}{c} \right)^{-1} U^T$$

s.t.  $K \succ 0$  SDP

$$\begin{aligned} \text{Tr}(K^P) &\leq B \\ -l \leq n_p - n_q &\leq l \\ L &= I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}} \end{aligned}$$

\left. \begin{array}{l} \text{avoid overfitting} \\ \text{and all points in the cluster} \end{array} \right\}

(Logistic)

This maintains the spectrum of the data after regularization.

$O(n^2)$  to  $O(n^3)$  complexity.

### Refs

#### Kernel Theory

- [1] Statistical Learning Theory, V.N. Vapnik
- [2] Learning with Kernels, B. Scholkopf and A. Smola

#### Matrix Approximation

- [3] On the Nyström Method for Approximating a Gram Matrix for Improved Kernel-Based Learning, Drineas and Mahoney
- [4] CUR matrix decompositions for improved data analysis, Drineas and Mahoney
- [5] Improving CUR Matrix Decomposition and the Nyström Approximation via Adaptive Sampling, Wang and Zhang
- [6] Fast Computation of Low Rank Matrix Approximations, Achlioptas and McSherry

#### Random Projections

- [7] Kernels as Features: On kernels, margins, and low-dimensional mappings, Balcan, Blum and Vempala
- [8] Random Features for Large-Scale Kernel Machines, Rahimi and Recht

#### Kernel Learning

- [9] Generalized Maximum Margin Clustering and Kernel Learning, Valizadegan and Jin