Data Analysis and Manifold Learning Lecture 8: A Brief Introduction to Kernel Methods

> Radu Horaud INRIA Grenoble Rhone-Alpes, France Radu.Horaud@inrialpes.fr http://perception.inrialpes.fr/

Outline of Lecture 8

- Linear regression in "feature space"
- Kernel construction and characterization of the feature space.
- The kernel (Gram) matrix
- The covariance matrix in feature space
- Feature-space computations
- Kernal PCA

Material for this lecture

- C. Bishop. Pattern Analysis and Machine Learning (chapters 6 and 12).
- J. Shawe-Taylor & N. Cristianini. Kernel Methods in Pattern Analysis (chapters 2, 3, 5 and 6).

The Kernel Function

- Consider a data set: $\mathbf{X} = [oldsymbol{x}_1, \dots, oldsymbol{x}_n] \in \mathbb{R}^D$
- Definition of a kernel function: consider a nonlinear feature space mapping: $\phi : x \to \phi(x)$, with $\phi(x) \in \mathbb{R}^M$. A kernel satisfies:

$$\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \boldsymbol{\phi}(\boldsymbol{x}_i)^{\top} \boldsymbol{\phi}(\boldsymbol{x}_j) = \langle \boldsymbol{\phi}(\boldsymbol{x}_i), \boldsymbol{\phi}(\boldsymbol{x}_j)
angle$$

- The main principle of *kernel methods* is to interpret the kernel function as an *inner product* in feature space and to design algorithms without making explicit the function φ.
- This extends many algorithms by making use of the *kernel trick* or *kernel substitution*.
- For example, we can extend basic algorithms, such as PCA and LDA in feature space, namely *kernel PCA* and *kernel Fisher discriminant*, etc.

Linear Regression in Feature Space

• Replace the standard regression problem with:

$$y = \sum_{m=1}^{M} w_m \phi_m(\boldsymbol{x}) = \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x})$$

• The parameters w_1, \ldots, w_M can be estimated from a training set of pairs (y_j, x_j) by minimizing the following criterion:

$$J(\boldsymbol{w}) = \frac{1}{2} \left(\sum_{j=1}^{n} (\boldsymbol{w}^{\top} \boldsymbol{\phi}(\boldsymbol{x}_{j}) - y_{j})^{2} + \lambda \boldsymbol{w}^{\top} \boldsymbol{w} \right)$$

Least-square Solution

• By taking the derivatives of J with respect to w and setting them to zero, we obtain the following solution:

$$oldsymbol{w} = -rac{1}{\lambda}\sum_{j=1}^n (oldsymbol{w}^ op oldsymbol{\phi}(oldsymbol{x}_j) - y_j) oldsymbol{\phi}(oldsymbol{x}_j)$$

- Let $a_j = -\frac{1}{\lambda} (\boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x}_j) y_j)$ be the *j*-th entry of a vector $\boldsymbol{a} \in \mathbb{R}^n$.
- Let $\Phi = [\phi(x_1) \dots \phi(x_j) \dots \phi(x_n)]$ be a $M \times n$ data matrix in feature space.
- Hence:

$$w = \Phi a$$

• We will use *a* instead of *w*.

Dual representation

• Substitute $w = \Phi a$ in J(w). We obtain:

$$J(\boldsymbol{a}) = \frac{1}{2} \left(\boldsymbol{a}^{\top} \boldsymbol{\Phi}^{\top} \boldsymbol{\Phi} \boldsymbol{\Phi}^{\top} \boldsymbol{\Phi} \boldsymbol{a} - \boldsymbol{a}^{\top} \boldsymbol{\Phi}^{\top} \boldsymbol{\Phi} \boldsymbol{y} + \lambda \boldsymbol{a}^{\top} \boldsymbol{\Phi}^{\top} \boldsymbol{\Phi} \boldsymbol{a} \right)$$

$$\mathbf{K} = \mathbf{\Phi}^{ op} \mathbf{\Phi}$$

is a Gram matrix in feature space (it will be referred to as a kernel matrix), with entries:

$$\kappa_{ij} = \langle \boldsymbol{\phi}(\boldsymbol{x}_i), \boldsymbol{\phi}(\boldsymbol{x}_j) \rangle = \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

Solution in feature space

 We obtain a new expression for J(w) as a function of vector a and of the Gram matrix:

$$J(\boldsymbol{a}) = \frac{1}{2} \left(\boldsymbol{a}^{\top} \mathbf{K}^{\top} \mathbf{K} \boldsymbol{a} - \boldsymbol{a}^{\top} \mathbf{K} \boldsymbol{y} + \lambda \boldsymbol{a}^{\top} \mathbf{K} \boldsymbol{a} \right)$$

• The solution is obtained by setting the gradient of J with respect to a to zero:

$$\boldsymbol{a} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \boldsymbol{y}$$

which always has an inverse.

Back to linear regression

• The linear regression model allows a prediction for a new input *x*:

$$y = \sum_{m=1}^{M} w_m \phi_m(\boldsymbol{x}) = \boldsymbol{w}^\top \boldsymbol{\phi}(\boldsymbol{x})$$

• By substitution this becomes:

$$y = \boldsymbol{a}^{\top} \boldsymbol{\Phi}^{\top} \boldsymbol{\phi}(\boldsymbol{x}) = \sum_{j=1}^{n} a_j \kappa(\boldsymbol{x}_j, \boldsymbol{x})$$

Discussion

- The dual representation allows the solution to be expressed entirely in terms of the kernel function;
- Inversion of a $M \times M$ matrix (dimension of the feature space) is replaced by inversion of a $n \times n$ matrix (number of points in the training set).
- $\bullet\,$ It avoids computations in feature space when M is very large.
- The feature-space is a vector space equipped with an inner-product – metric space;
- This means that there is a strong similarity between feature-space methods and MDS (only the pairwise inner-product between data points are needed to construct algorithms).

Constructing Kernels

- A valid kernel function is such that the associated Gram matrix is symmetric positive semidefinite.
- ullet The simplest kernel corresponds to $\phi(m{x})=m{x}$, or

$$\kappa(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{x}^\top \boldsymbol{x}'$$

Valid kernels:

$$\begin{split} c\kappa(\boldsymbol{x},\boldsymbol{x}') \text{ with } c > 0; \quad & f(\boldsymbol{x})\kappa(\boldsymbol{x},\boldsymbol{x}')f(\boldsymbol{x}'); \quad \exp(\kappa(\boldsymbol{x},\boldsymbol{x}')); \\ \kappa_1(\boldsymbol{x},\boldsymbol{x}') + \kappa_2(\boldsymbol{x},\boldsymbol{x}'); \quad & \kappa_1(\boldsymbol{x},\boldsymbol{x}')\kappa_2(\boldsymbol{x},\boldsymbol{x}'); \\ \kappa(\boldsymbol{\phi}(\boldsymbol{x}),\boldsymbol{\phi}(\boldsymbol{x}')); \quad & \boldsymbol{x}^\top \mathbf{A} \boldsymbol{x}' \text{ with } \mathbf{A} \succeq 0. \end{split}$$

Kernel Normalization

• $oldsymbol{x}
ightarrow \phi(oldsymbol{x}) / \| oldsymbol{\phi}(oldsymbol{x}\|$ which yields:

$$\begin{split} \hat{\kappa}(\boldsymbol{x}, \boldsymbol{x}') &= \frac{\kappa(\boldsymbol{x}, \boldsymbol{x}')}{\sqrt{\kappa(\boldsymbol{x}, \boldsymbol{x})\kappa(\boldsymbol{x}', \boldsymbol{x}')}} \\ &= \kappa(\boldsymbol{x}, \boldsymbol{x})^{-1/2}\kappa(\boldsymbol{x}, \boldsymbol{x}')\kappa(\boldsymbol{x}', \boldsymbol{x}')^{-1/2} \end{split}$$

The Gaussian Kernel

$$\kappa(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-rac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2\sigma^2}
ight)$$

•
$$\|x - x'\|^2 = x^\top x + x'^\top x' - 2x^\top x'$$

• Let $f(x) = \exp(-x^\top x/2\sigma^2) = \frac{1}{\sqrt{\exp(x^\top x/\sigma^2)}}$

• The Gaussian kernel writes:

$$\begin{split} \kappa(\boldsymbol{x}, \boldsymbol{x}') &= f(\boldsymbol{x}) \exp(\boldsymbol{x}^{\top} \boldsymbol{x}' / \sigma^2) f(\boldsymbol{x}') \\ &= \frac{\exp(\boldsymbol{x}^{\top} \boldsymbol{x}' / \sigma^2)}{\sqrt{\exp(\boldsymbol{x}^{\top} \boldsymbol{x} / \sigma^2) \exp(\boldsymbol{x}'^{\top} \boldsymbol{x}' / \sigma^2)}} \end{split}$$

• This is also known as the basis radial function (BRF) kernel.

Mercer Kernel (In Brief!)

• Let X be a compact subset of \mathbb{R}^D . Suppose that k is a continuous and symmetric function such that the integral operator is positive

$$\int_{X \times X} \kappa(\boldsymbol{x}, \boldsymbol{x}') f(\boldsymbol{x}) f(\boldsymbol{x}') d\boldsymbol{x} d\boldsymbol{x}' \geq 0$$

for all $f \in L_2(X)$. An L_2 function is a function that is square integrable.

• We can expand $\kappa(\boldsymbol{x}, \boldsymbol{x}')$ in a uniformly convergent series in terms of functions $\{\phi_i\}_{i=1}^{\infty}$ satisfying $\langle \phi_i, \phi_j \rangle = \delta_{ij}$

$$\kappa({m x},{m x}') = \sum_{i=1}^\infty \phi_i({m x}) \phi_i({m x}')$$

Inner-product Space

• A vector space is an inner-product space if there exists a real-valued symmetric bilinear map that satisfies:

$$\langle \boldsymbol{x}, \boldsymbol{x} \rangle \geq 0$$

- The inner product is *strict* if: $\langle \boldsymbol{x}, \boldsymbol{x} \rangle = 0$ iff $\boldsymbol{x} = 0$.
- A strict inner product allows to define a norm of a vector $\|x\|_2 = \sqrt{\langle x, x \rangle}$ and an associated *metric* or distance $\|x x'\|_2$.
- A vector space with a metric is known as a metric space.
- The feature space is a metric space, equipped with the strict inner product.

The Gram/Kernel Matrix

• The $n \times n$ matrix:

$$\mathbf{K} = \mathbf{\Phi}^\top \mathbf{\Phi}$$

is a Gram matrix in feature space, with entries:

$$\kappa_{ij} = \langle \boldsymbol{\phi}(\boldsymbol{x}_i), \boldsymbol{\phi}(\boldsymbol{x}_j) \rangle = \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

- We remind that $\Phi = [\phi(x_1) \dots \phi(x_j) \dots \phi(x_n)]$ is the $M \times n$ data matrix.
- It is symmetric, positive, semidefinite:

$$oldsymbol{x}^{ op}\mathbf{K}oldsymbol{x}=oldsymbol{x}^{ op}\mathbf{\Phi}oldsymbol{x}=\|oldsymbol{\Phi}oldsymbol{x}\|_2^2$$

 This matrix was studied in Lecture #1 within the context of MDS. Here we have a generalization because each entry is a kernel function which is more general hat the dot-product of MDS.

Spectral Decomposition of the Kernel Matrix

 Let (λ₁, v₁),..., (λ_n, v_n) be the eigenvalue-eigenvector pairs of a Kernel matrix. It can be written as:

$$\mathbf{K} = \sum_{k=1}^n \lambda_k oldsymbol{v}_k oldsymbol{v}_k^ op$$

• Each matrix entry can be written as:

$$\kappa_{ij} = \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \sum_{k=1}^n \lambda_k v_{ik} v_{jk} = \langle \boldsymbol{\phi}(\boldsymbol{x}_i), \boldsymbol{\phi}(\boldsymbol{x}_j) \rangle$$

- with $\boldsymbol{\phi}(\boldsymbol{x}_i) = (\sqrt{\lambda_1} v_{i1}, \dots, \sqrt{\lambda_k} v_{ik}, \dots, \sqrt{\lambda_n} v_{in})^\top$.
- Therefore, we can think of the eigenvectors as defining a feature space.

Feature-space Computations

• The norm of a feature-space vector:

$$\|oldsymbol{\phi}(oldsymbol{x})\|_2^2 = \langle oldsymbol{\phi}(oldsymbol{x}), oldsymbol{\phi}(oldsymbol{x})
angle = \kappa(oldsymbol{x},oldsymbol{x})$$

• The norm of a linear combination:

$$\|\sum_{i=1}^{n} \alpha_i \boldsymbol{\phi}(\boldsymbol{x}_i)\|^2 = \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

• Distance between two feature-space vectors:

$$\| \boldsymbol{\phi}(\boldsymbol{x}_i) - \boldsymbol{\phi}(\boldsymbol{x}_j) \|_2^2 = \kappa(\boldsymbol{x}_i, \boldsymbol{x}_i) + \kappa(\boldsymbol{x}_j, \boldsymbol{x}_j) - 2\kappa(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

Center of Mass

 $\frac{1}{n}$

• Notation:
$$\overline{oldsymbol{\phi}} = rac{1}{n} \sum_{i=1}^n oldsymbol{\phi}(oldsymbol{x}_i)$$

- There is no explicit dual representation for this point. Moreover, it is not the image of a "valid" data point.
- Norm, distance from a point, and expected distance:

$$\begin{split} \|\overline{\boldsymbol{\phi}}\|^2 &= \sum_{i=1}^n \sum_{j=1}^n \frac{1}{n} \frac{1}{n} \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) \\ \|\boldsymbol{\phi}(\boldsymbol{x}) - \overline{\boldsymbol{\phi}}\|^2 &= \langle \boldsymbol{\phi}(\boldsymbol{x}), \boldsymbol{\phi}(\boldsymbol{x}) \rangle + \langle \overline{\boldsymbol{\phi}}, \overline{\boldsymbol{\phi}} \rangle - 2 \langle \boldsymbol{\phi}(\boldsymbol{x}), \overline{\boldsymbol{\phi}} \rangle \\ &= \kappa(\boldsymbol{x}, \boldsymbol{x}) + \frac{1}{n^2} \sum_{i,j=1}^n \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) - \frac{2}{n} \sum_{i=1}^n \kappa(\boldsymbol{x}, \boldsymbol{x}_i) \\ \sum_{k=1}^n \|\boldsymbol{\phi}(\boldsymbol{x}_k) - \overline{\boldsymbol{\phi}}\|^2 &= \frac{1}{n} \sum_{k=1}^n \kappa(\boldsymbol{x}_k, \boldsymbol{x}_k) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \kappa(\boldsymbol{x}_i, \boldsymbol{x}_j) \end{split}$$

The Kernel Matrix of Centered Data

• In feature-space the centered data writes:

$$\hat{oldsymbol{\phi}}(oldsymbol{x}) = oldsymbol{\phi}(oldsymbol{x}) - \overline{oldsymbol{\phi}}$$

• The corresponding entry of the associated kernel matrix writes:

$$\begin{split} \hat{\kappa}(\boldsymbol{x},\boldsymbol{x}') &= \langle \boldsymbol{\phi}(\boldsymbol{x}) - \overline{\boldsymbol{\phi}}, \boldsymbol{\phi}(\boldsymbol{x}') - \overline{\boldsymbol{\phi}} \rangle \\ &= \kappa(\boldsymbol{x},\boldsymbol{x}') - \frac{1}{n} \sum_{i=1}^{n} \left(\kappa(\boldsymbol{x},\boldsymbol{x}_i) - \kappa(\boldsymbol{x}',\boldsymbol{x}_i) \right) + \frac{1}{n^2} \sum_{i,j=1}^{n} \kappa(\boldsymbol{x}_i,\boldsymbol{x}_j) \end{split}$$

• In matrix form:

$$\widehat{\mathbf{K}} = \mathbf{K} - \frac{1}{n} \sum_{i=1}^{n} \left(\mathbb{1}\mathbb{1}^{\top}\mathbf{K} + \mathbf{K}\mathbb{1}\mathbb{1}^{\top} \right) + \frac{1}{n^2} (\mathbb{1}^{\top}\mathbf{K}\mathbb{1})\mathbb{1}\mathbb{1}^{\top}$$

The Spread of the Data

• The $M \times n$ data matrix in feature space:

$$oldsymbol{\Phi} = [oldsymbol{\phi}(oldsymbol{x}_1)\dotsoldsymbol{\phi}(oldsymbol{x}_n)]$$

• The covariance matrix for centered data is an $M \times M$ matrix:

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

• Each entry of this matrix is:

$$c_{st} = rac{1}{n}\sum_{i=1}^n oldsymbol{\phi}(oldsymbol{x}_i)_s \ oldsymbol{\phi}(oldsymbol{x}_i)_t$$

The Projected Variance

• For centered data, the variance along a vector v writes:

$$\sigma_{\boldsymbol{v}}^2 = \frac{1}{n} \boldsymbol{v}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{v}$$

• If the data are not centered:

$$\sigma_{\boldsymbol{v}}^2 = \frac{1}{n} \boldsymbol{v}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{v} - \left(\frac{1}{n} \boldsymbol{v}^\top \boldsymbol{\Phi} \mathbb{1}\right)^2$$

Dual Representation of the Projected Variance

- Let's write v as a combination of the feature-space points: $v = \sum_{i=1}^{n} \alpha_i \phi(x_i) = \Phi \alpha.$
- By substitution in the formula of the projected variance, we obtain:

$$\begin{aligned} \sigma_{\boldsymbol{v}}^2 &= \frac{1}{n} \boldsymbol{\alpha}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \boldsymbol{\alpha} - \left(\frac{1}{n} \boldsymbol{\alpha}^\top \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \mathbb{1}\right)^2 \\ &= \frac{1}{n} \boldsymbol{\alpha}^\top \mathbf{K}^2 \boldsymbol{\alpha} - \left(\frac{1}{n} \boldsymbol{\alpha}^\top \mathbf{K} \mathbb{1}\right)^2 \end{aligned}$$

Eigendecomposition of Covariance and Kernel Matrices

• For centred data we have:

$$\begin{split} \mathbf{C} &= \quad \frac{1}{n} \mathbf{\Phi} \mathbf{\Phi}^\top \text{ with } \{(\mu_i, \boldsymbol{u}_i)\}_{i=1}^M \\ \mathbf{K} &= \quad \mathbf{\Phi}^\top \mathbf{\Phi} \text{ with } \{(\lambda_i, \boldsymbol{v}_i)\}_{i=1}^n \end{split}$$

• By premultiplication of $\mathbf{\Phi}\mathbf{\Phi}^{\top}\mathbf{u} = n\mu\mathbf{u}$ with $\mathbf{\Phi}^{\top}$ we obtain:

$$oldsymbol{v} = oldsymbol{\Phi}^{ op} oldsymbol{u}$$
 and $oldsymbol{\lambda} = n \mu$

- From which we obtain: $\| \boldsymbol{v} \|^2 = \boldsymbol{u}^\top \boldsymbol{\Phi} \boldsymbol{\Phi}^\top \boldsymbol{u} = n \mu = \lambda$
- The normalized eigenvector of the kernel matrix is:

$$oldsymbol{v} = \lambda^{-1/2} oldsymbol{\Phi}^{ op} oldsymbol{u}$$

• There is a similar dual expression:

$$\boldsymbol{u} = \lambda^{-1/2} \boldsymbol{\Phi} \boldsymbol{v}$$

Traces

• The traces are related by:

$$\mathsf{tr}(\mathbf{C}) = rac{1}{n}\mathsf{tr}(\mathbf{K})$$

• The trace of the kernel matrix:

$$\mathsf{tr}(\mathbf{K}) = \sum_{i=1}^n \kappa(oldsymbol{x}_i, oldsymbol{x}_i)$$

• The total variance in feature-space:

$$\sum_{i=1}^{M} \mu_i = \frac{1}{n} \sum_{i=1}^{n} \kappa(\boldsymbol{x}_i, \boldsymbol{x}_i)$$

• This can be used to estimate the dimension $m \ll M$ of the reduced feature space.

Covariance Eigenvectors in Feature-space

• The eigenvectors of the covariance matrix:

$$\mathbf{U} = \left[\begin{array}{ccc} \lambda_1^{-1/2} \mathbf{\Phi} \boldsymbol{v}_1 & \dots & \lambda_k^{-1/2} \mathbf{\Phi} \boldsymbol{v}_k & \dots \end{array}\right] = \mathbf{\Phi} \mathbf{V} \boldsymbol{\Lambda}^{-1/2}$$

• Each eigenvector:

$$\boldsymbol{u}_k = \lambda_k^{-1/2} \boldsymbol{\Phi} \boldsymbol{v}_k = \lambda_k^{-1/2} \sum_{i=1}^n v_{ik} \boldsymbol{\phi}(\boldsymbol{x}_i)$$

• Let:
$$\boldsymbol{\beta}_k = \lambda_k^{-1/2} \boldsymbol{v}_k = (\lambda_k^{-1/2} v_{1k} \dots \lambda_k^{-1/2} v_{ik} \dots \lambda_k^{-1/2} v_{nk})$$

Hence:

$$oldsymbol{u}_k = \sum_{i=1}^n eta_{ik} oldsymbol{\phi}(oldsymbol{x}_i)$$

Projection of a data point on a principal direction

• Let's project a data point in feature space $\phi(x)$ onto an eigenvector of the covariance matrix:

$$egin{aligned} oldsymbol{u}_k^ op oldsymbol{\phi}(oldsymbol{x}) &=& \langleoldsymbol{u}_k, oldsymbol{\phi}(oldsymbol{x})
angle \ &=& \sum_{i=1}^n eta_{ik} \langleoldsymbol{\phi}(oldsymbol{x}_i), oldsymbol{\phi}(oldsymbol{x})
angle \ &=& \sum_{i=1}^n eta_{ik} \kappa(oldsymbol{x}_i,oldsymbol{x}) \end{aligned}$$

 Let U be the M × m matrix formed with m ≪ M eigenvectors of C. A feature point can be mapped in the eigenspace of C with:

$$ilde{oldsymbol{\phi}}(oldsymbol{x}) = \mathbf{U}^{ op} oldsymbol{\phi}(oldsymbol{x})$$

Kernel PCA

• Consider the centered kernel matrix:

$$\widehat{\mathbf{K}} = \mathbf{K} - \frac{1}{n} \sum_{i=1}^{n} \left(\mathbb{1}\mathbb{1}^{\top} \mathbf{K} + \mathbf{K}\mathbb{1}\mathbb{1}^{\top} \right) + \frac{1}{n^2} (\mathbb{1}^{\top} \mathbf{K}\mathbb{1})\mathbb{1}\mathbb{1}^{\top}$$

• Compute the eigen decomposition of this matrix and retain the *m* largest eigenvalue-eigenvector pairs:

$$\widehat{\mathbf{K}} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^{\top}$$

- Compute the vectors: $oldsymbol{eta}_k = \lambda_k^{-1/2} oldsymbol{v}_k, \; k=1\dots m$
- Transform the data:

$$ilde{oldsymbol{\phi}}(oldsymbol{x}) = \mathbf{U}^{ op} oldsymbol{\phi}(oldsymbol{x})$$

Additional Topics of Interest

- Kernel K-means clustering http://citeseerx.ist.psu.edu/viewdoc/download? doi=10.1.1.79.2967&rep=rep1&type=pdf
- Kernel Fisher discriminant analysis http://ieeexplore.ieee.org/xpls/abs_all.jsp? arnumber=788121
- Diffusion kernels (next course)