

Manifold Learning for Signal and Image Analysis

Lecture 5: A Brief Introduction to Kernel Methods

Radu Horaud

INRIA Grenoble Rhone-Alpes, France

Radu.Horaud@inria.fr

<http://perception.inrialpes.fr/>

Outline of Lecture 5

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

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

- 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:

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

- The parameters w_1, \dots, w_M can be estimated from a training set of pairs (y_j, \mathbf{x}_j) by minimizing the following criterion:

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

Least-square Solution

- By taking the derivatives of J with respect to \mathbf{w} and setting them to zero, we obtain the following solution:

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{j=1}^n (\mathbf{w}^\top \phi(\mathbf{x}_j) - y_j) \phi(\mathbf{x}_j)$$

- Let $a_j = -\frac{1}{\lambda} (\mathbf{w}^\top \phi(\mathbf{x}_j) - y_j)$ be the j -th entry of a vector $\mathbf{a} \in \mathbb{R}^n$.
- Let $\Phi = [\phi(\mathbf{x}_1) \dots \phi(\mathbf{x}_j) \dots \phi(\mathbf{x}_n)]$ be a $M \times n$ data matrix in feature space.
- Hence:

$$\mathbf{w} = \Phi \mathbf{a}$$

- We will use \mathbf{a} instead of \mathbf{w} .

Dual representation

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

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

- The $n \times n$ matrix:

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

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

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

Solution in Feature Space

- We obtain a new expression for $J(\mathbf{w})$ as a function of vector \mathbf{a} and of the Gram matrix:

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

- The solution is obtained by setting the gradient of J with respect to \mathbf{a} to zero:

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

which always has an inverse.

Back to Linear Regression

- The linear regression model allows to predict the output y from a new input \mathbf{x} :

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

- By substitution this becomes:

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

- The mapping ϕ is not required!

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).
- 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.
- The simplest kernel corresponds to $\phi(\mathbf{x}) = \mathbf{x}$, or

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

- Valid kernels:

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

Kernel Normalization

- $\mathbf{x} \rightarrow \phi(\mathbf{x})/\|\phi(\mathbf{x})\|$ which yields:

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

The Gaussian Kernel

$$\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$$

- $\|\mathbf{x} - \mathbf{x}'\|^2 = \mathbf{x}^\top \mathbf{x} + \mathbf{x}'^\top \mathbf{x}' - 2\mathbf{x}^\top \mathbf{x}'$
- Let $f(\mathbf{x}) = \exp(-\mathbf{x}^\top \mathbf{x}/2\sigma^2) = \frac{1}{\sqrt{\exp(\mathbf{x}^\top \mathbf{x}/\sigma^2)}}$
- The Gaussian kernel writes:

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

- 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

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

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

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

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

Inner-product Space

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

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

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

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

- We remind that $\mathbf{\Phi} = [\phi(\mathbf{x}_1) \dots \phi(\mathbf{x}_j) \dots \phi(\mathbf{x}_n)]$ is the $M \times n$ data matrix.
- It is symmetric, positive, semidefinite:

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

- 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 than the dot-product of MDS.

Spectral Decomposition of the Kernel Matrix

- Let $(\lambda_1, \mathbf{v}_1), \dots, (\lambda_n, \mathbf{v}_n)$ be the eigenvalue-eigenvector pairs of a Kernel matrix. It can be written as:

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

- Each matrix entry can be written as:

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

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

Feature-space Computations

- The norm of a feature-space vector:

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

- The norm of a linear combination:

$$\left\| \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) \right\|_2^2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

- Distance between two feature-space vectors:

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

Center of Mass

- Notation: $\bar{\phi} = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i)$
- 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:

$$\|\bar{\phi}\|^2 = \sum_{i=1}^n \sum_{j=1}^n \frac{1}{n} \frac{1}{n} \kappa(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

$$\begin{aligned} \|\phi(\mathbf{x}) - \bar{\phi}\|^2 &= \langle \phi(\mathbf{x}), \phi(\mathbf{x}) \rangle + \langle \bar{\phi}, \bar{\phi} \rangle - 2\langle \phi(\mathbf{x}), \bar{\phi} \rangle \\ &= \kappa(\mathbf{x}, \mathbf{x}) + \frac{1}{n^2} \sum_{i,j=1}^n \kappa(\mathbf{x}_i, \mathbf{x}_j) - \frac{2}{n} \sum_{i=1}^n \kappa(\mathbf{x}, \mathbf{x}_i) \end{aligned}$$

$$\frac{1}{n} \sum_{k=1}^n \|\phi(\mathbf{x}_k) - \bar{\phi}\|^2 = \frac{1}{n} \sum_{k=1}^n \kappa(\mathbf{x}_k, \mathbf{x}_k) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \kappa(\mathbf{x}_i, \mathbf{x}_j)$$

The Kernel Matrix of Centered Data

- In feature-space the centered data writes:

$$\hat{\phi}(\mathbf{x}) = \phi(\mathbf{x}) - \bar{\phi}$$

- The corresponding entry of the associated kernel matrix writes:

$$\begin{aligned}\hat{\kappa}(\mathbf{x}, \mathbf{x}') &= \langle \phi(\mathbf{x}) - \bar{\phi}, \phi(\mathbf{x}') - \bar{\phi} \rangle \\ &= \kappa(\mathbf{x}, \mathbf{x}') - \frac{1}{n} \sum_{i=1}^n (\kappa(\mathbf{x}, \mathbf{x}_i) + \kappa(\mathbf{x}', \mathbf{x}_i)) + \frac{1}{n^2} \sum_{i,j=1}^n \kappa(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}$$

- In matrix form:

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

The Spread of the Data

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

$$\Phi = [\phi(\mathbf{x}_1) \dots \phi(\mathbf{x}_j) \dots \phi(\mathbf{x}_n)]$$

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

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

- Each entry of this matrix is:

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

The Projected Variance

- For centered data, the variance along a vector \mathbf{v} writes:

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

- If the data are not centered:

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

Dual Representation of the Projected Variance

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

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

Eigendecomposition of Covariance and Kernel Matrices

- For centred data we have:

$$\mathbf{C} = \frac{1}{n} \mathbf{\Phi} \mathbf{\Phi}^\top \text{ with } \{(\mu_i, \mathbf{u}_i)\}_{i=1}^M$$

$$\mathbf{K} = \mathbf{\Phi}^\top \mathbf{\Phi} \text{ with } \{(\lambda_i, \mathbf{v}_i)\}_{i=1}^n$$

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

$$\mathbf{v} = \mathbf{\Phi}^\top \mathbf{u} \text{ and } \lambda = n\mu$$

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

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

- There is a similar dual expression:

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

Traces

- The traces are related by:

$$\text{tr}(\mathbf{C}) = \frac{1}{n} \text{tr}(\mathbf{K})$$

- The trace of the kernel matrix:

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

- The total variance in feature-space:

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

- 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:

$$\mathbf{U} = \left[\lambda_1^{-1/2} \Phi \mathbf{v}_1 \quad \dots \quad \lambda_k^{-1/2} \Phi \mathbf{v}_k \quad \dots \right] = \Phi \mathbf{V} \Lambda^{-1/2}$$

- Each eigenvector:

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

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

- Hence:

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

Projection of a Data Point on a Principal Direction

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

$$\begin{aligned}\mathbf{u}_k^\top \phi(\mathbf{x}) &= \langle \mathbf{u}_k, \phi(\mathbf{x}) \rangle \\ &= \sum_{i=1}^n \beta_{ik} \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle \\ &= \sum_{i=1}^n \beta_{ik} \kappa(\mathbf{x}_i, \mathbf{x})\end{aligned}$$

- Let \mathbf{U} be the $M \times m$ matrix formed with $m \ll M$ eigenvectors of \mathbf{C} . A feature point can be mapped in the eigenspace of \mathbf{C} with:

$$\tilde{\phi}(\mathbf{x}) = \mathbf{U}^\top \phi(\mathbf{x})$$

The Kernel PCA method

- Build the centered kernel matrix associated with a data set \mathbf{X} and a kernel $\kappa(\mathbf{x}, \mathbf{x}')$:

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

- Compute the eigen-decomposition of this matrix and retain the K largest eigenvalue-eigenvector pairs:

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

- Compute the vectors: $\beta_k = \lambda_k^{-1/2} \mathbf{v}_k$, $k = 1 \dots K$
- Project the feature-space data onto the space spanned by the eigenvectors of the feature-space covariance:

$$\tilde{\phi}(\mathbf{x}) = \mathbf{U}^\top \phi(\mathbf{x}) \quad \text{with} \quad \mathbf{u}_k^\top \phi(\mathbf{x}) = \sum_{i=1}^n \beta_{ik} \kappa(\mathbf{x}_i, \mathbf{x})$$

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 and exponential kernels (next lecture)