

Lecture 7 — February 11

Lecturer: Martin Wainwright

Scribe: Vivek Ramamurthy

Note: These lecture notes are still rough, and have only have been mildly proofread.



This is the danger environment.

7.1 Announcements

HW #2: due Monday February 23.

7.2 Outline

- Mercer's characterization
- Kernel PCA (dimensionality reduction)

7.3 Mercer's characterization

Given a symmetric and positive semidefinite matrix $K \in \mathbb{R}^{d \times d}$, we know from standard linear algebra that there exist real scalars $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ and vectors $\{\psi_i, i = 1, \dots, d\}$ such that

$$K = \sum_{i=1}^d \lambda_i \psi_i \psi_i^T$$

In this decomposition, the vectors $\{\psi_i\}$ are eigenvectors, obtained by solving the matrix-vector equation

$$K\psi = \lambda\psi.$$

Moreover, the $\{\psi_i\}$ can be chosen to be an orthonormal system of vectors.

We now discuss a generalization of this type of decomposition to the more general setting of linear operators in a Hilbert space. (The matrix is a special case of a linear operator on \mathbb{R}^d .)

Given a Hilbert space $\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$ of functions, a linear operator $T : \mathcal{H} \rightarrow \mathcal{H}$ is a mapping such that

1. $\forall f \in \mathcal{H}, T(f) \in \mathcal{H}$
2. $\forall f, g \in \mathcal{H}, T(f + g) = T(f) + T(g)$
3. $\forall \alpha \in \mathbb{R}, T(\alpha f) = \alpha T(f)$

7.3.1 Mercer's theorem (one variant):

Theorem 7.1. Say $\mathcal{X} \subseteq \mathbb{R}^d$ is compact, and $\mathbb{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is continuous, and satisfies

$$\int_y \int_x \mathbb{K}^2(x, y) dx dy < +\infty,$$

$$\int_y \int_x f(x) \mathbb{K}(x, y) f(y) dx dy \geq 0 \quad \forall f \in L^2(\mathcal{X}) \text{ (i.e. a positive semidefinite kernel)}$$

$$\text{where } L^2(\mathcal{X}) = \{f : \int f^2(x) dx < +\infty\}$$

Then there exist $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots$ (all non-negative) and functions $\{\psi_i(\cdot) \in L^2(\mathcal{X}), i = 1, 2, 3, \dots\}$ such that

$$\mathbb{K}(x, y) = \sum_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(y) \quad \forall x, y \in \mathcal{X}$$

Moreover, the $\{\psi_i\}$ are an orthonormal system in $L^2(\mathcal{X})$, meaning that

$$\langle \psi_i, \psi_j \rangle_{L^2(\mathcal{X})} = \int \psi_i(x) \psi_j(x) dx = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

Remarks: Note that this can be seen as a generalization of the decomposition

$$\mathbb{K}(x, y) = x^T K y = \sum_{i=1}^d \lambda_i (\psi_i^T x) (\psi_i^T y)$$

in the finite-dimensional setting. The orthogonality condition is a generalization of the fact that PSD matrices have an orthogonal set of eigenvectors.

Mercer's theorem is a special case of spectral decomposition theory for self-adjoint, positive operators in Hilbert spaces.

7.3.2 Use of Mercer's Theorem

Eigenfunctions can be obtained by solving the integral equation:

$$T_{\mathbb{K}}(f)(x) := \int \mathbb{K}(x, y) f(y) dy = \lambda f(x)$$

Here

$$T_{\mathbb{K}}(f)(\cdot) := \int \mathbb{K}(\cdot, y) f(y) dy$$

is a linear operator on $L^2(\mathcal{X}) \rightarrow L^2(\mathcal{X})$. (Homework #2 has some instances of this procedure.)

We can then use the eigenfunctions thus obtained to generate a “feature map” given by

$$\Phi : \mathcal{X} \rightarrow \ell^2(\mathbb{N})$$

Here, the feature map Φ maps data $x \in \mathcal{X}$ to a sequence $(a_1, a_2, \dots) \in \ell^2(\mathbb{N})$, where

$$\ell^2(\mathbb{N}) = \{(a_1, a_2, \dots) \mid \sum_{i=1}^{\infty} a_i^2 < +\infty\}$$

For example, consider the feature map defined as follows:

$$\Phi(x) = (\sqrt{\lambda_1}\psi_1(x), \sqrt{\lambda_2}\psi_2(x), \dots, \sqrt{\lambda_i}\psi_i(x), \dots).$$

That is, we map each $x \in \mathcal{X}$ into a sequence $\Phi(x)$ in $\ell^2(\mathbb{N})$.

Using Mercer’s decomposition, if we take the inner product (in $\ell^2(\mathbb{N})$) between the two sequences $\Phi(x)$ and $\Phi(y)$, then we recover the kernel function

$$\langle \Phi(x), \Phi(y) \rangle_{\ell^2(\mathbb{N})} = \sum_{i=1}^{\infty} \sqrt{\lambda_i}\psi_i(x)\sqrt{\lambda_i}\psi_i(y) = \mathbb{K}(x, y).$$

7.4 Kernel PCA

7.4.1 Quick recap on classical PCA

Given data $X^{(1)}, \dots, X^{(n)} \subseteq \mathbb{R}^d$, we first compute the sample covariance or correlation matrix, given by

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n X^{(i)} [X^{(i)}]^T$$

Then, we compute the eigenvectors corresponding to the top $k \ll d$ eigenvalues (in value). Using these eigenvectors, we project data $\mathbf{X} \in \mathbb{R}^d$, a large space, into \mathbb{R}^k , a much smaller space. Thus, the primary motivation for PCA is achieving a large reduction in the dimensionality of the data.

To gain some intuition for PCA, consider an idealized “noisy subspace” generative model, given by

$$\mathbf{x} = \mathbf{V}\mathbf{z} + \mathbf{w}$$

where $\mathbf{V} \in \mathbb{R}^{d \times k}$ is fixed, $\mathbf{z} \in \mathbb{R}^k$ is random, and also $\mathbf{w} \in \mathbb{R}^d$ is random. Furthermore, we assume that

$$\mathbb{E}(\mathbf{z}) = 0, \text{ Cov}(\mathbf{z}) = \alpha^2 \mathbf{I}_{k \times k}$$

$$\mathbb{E}(\mathbf{w}) = 0, \text{ Cov}(\mathbf{w}) = \sigma^2 \mathbf{I}_{d \times d}$$

Finally, we assume that \mathbf{z} and \mathbf{w} are independent. This gives us

$$\text{Cov}(\mathbf{x}) = \Sigma = \alpha^2 \mathbf{V}\mathbf{V}^T + \sigma^2 \mathbf{I}_{d \times d}$$

Now, we may think of \mathbf{V} as having k orthogonal columns, i.e.,

$$\mathbf{V} = (V_1, \dots, V_k)$$

We also have that

$$\Sigma V_j = (\alpha^2 + \sigma^2) V_j$$

i.e., the eigenvectors corresponding to the top k eigenvalues are $\{V_1, \dots, V_k\}$. Moreover, for fixed d , we have that

$$\|\hat{\Sigma}_n - \Sigma\|_2 = \max_{\|u\|_2=1} \|(\hat{\Sigma}_n - \Sigma)u\|_2 \rightarrow 0 \text{ as } n \rightarrow +\infty$$

where $\|\cdot\|_2$ denotes the spectral radius (max. absolute value over all eigenvalues).

7.4.2 Kernel PCA (Scholkopf et. al., 1997)

We once again consider an idealized model, this time in feature space \mathcal{F} , which is given by

$$\Phi(\mathbf{x}) = \sum_{j=1}^k z_j \Phi_j + \mathbf{w} \quad (7.1)$$

where $\Phi_j \in \mathcal{F}$ for all $j = 1, \dots, k$ and is fixed, while $\mathbf{z} \in \mathbb{R}^k$ and $\mathbf{w} \in \mathcal{F}$ are both random.

Example: Suppose that we worked with the feature map defined by a polynomial kernel $\mathbb{K}(x, y) = (1 + \langle x, y \rangle)^m$ for $x \in \mathbb{R}^d$. In the special case $m = 2$ and $d = 2$, one feature map for this kernel is given by

$$\Phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2)$$

so that

$$\langle \Phi(x), \Phi(y) \rangle = 1 + 2x_1y_1 + 2x_2y_2 + 2x_1x_2y_1y_2 + x_1^2 + y_1^2 + x_2^2y_2^2 = (1 + x_1y_1 + x_2y_2)^2$$

One particular example of the model (7.1) would be

$$\begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_1x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix} = z_1 \Phi_1 + \mathbf{w}.$$

This would model the data as lying near to some quadratic surface, determined by the choice of $\Phi_1 \in \mathbb{R}^6$. \square

For simplicity, let us assume that the generating vectors are orthonormal

$$\langle \Phi_i, \Phi_j \rangle_{\mathcal{F}} = 0 \text{ if } i \neq j$$

Now let us define the covariance operator associated with the random element $\Phi(\mathbf{x})$. For each j , we use $\Phi_j \otimes \Phi_j$ to denote a linear operator on \mathcal{F} defined as follows: given some $f \in \mathcal{F}$, it outputs a new $(\Phi_j \otimes \Phi_j)(f) \in \mathcal{F}$, given by

$$(\Phi_j \otimes \Phi_j)(f) = \langle \Phi_j, f \rangle_{\mathcal{F}} \Phi_j.$$

With this definition, the covariance operator is given by

$$\text{Cov}[\Phi(\mathbf{x})] = \sum_{j=1}^k \text{Var}(z_j) (\Phi_j \otimes \Phi_j) + \mathbb{E}[\mathbf{w} \otimes \mathbf{w}]$$

Since it is a linear combination of linear operators, it is also a linear operator on \mathcal{F} .

In particular, for any $f \in \mathcal{F}$, this covariance operator outputs a new element of \mathcal{F} , given by

$$\text{Cov}[\Phi(\mathbf{x})](f) = \sum_{j=1}^k \text{Var}(z_j) \langle \Phi_j, f \rangle_{\mathcal{F}} \Phi_j + \mathbb{E}[\mathbf{w} \otimes \mathbf{w}](f)$$

At this point, the intuition underlying KPCA is the same as the intuition underlying PCA. That is, *if we knew* the functions Φ_j , then given a new sample, we could:

- map it to the feature space via $\mathbf{x} \mapsto \Phi(\mathbf{x})$
- compute its co-ordinates in the linear span of $\{\Phi_j\}$ by computing the projections $\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}) \rangle_{\mathcal{F}}$ for $j = 1, \dots, k$.

In practice, we don't know the $\{\Phi_j\}$, but as with ordinary PCA, we can try to estimate them from data. Given samples $\mathbf{x}^{(i)}, i = 1, 2, \dots, n$, we can form the empirical covariance operator

$$\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}^{(i)}) \otimes \Phi(\mathbf{x}^{(i)})$$

We would like to find eigenfunctions $\hat{\Phi}$ such that

$$(\hat{\Sigma}_n)(\hat{\Phi}) = \lambda \hat{\Phi} \tag{7.2}$$

The question now is, how do we express the above equation in terms of kernels, i.e. how do we "kernelize" it? Towards this end, we make the following claim:

Claim: Any solution to (7.2) is of the form

$$\hat{\Phi} = \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}^{(i)})$$

for some weight vector $(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$.

Proof: First, we observe that any solution to (7.2) lies in $\text{Range}(\widehat{\Sigma}_n)$. Linearity, and the nature of $\Phi(\mathbf{x}^{(i)}) \otimes \Phi(\mathbf{x}^{(i)})$ tell us that

$$\widehat{\Sigma}_n(\widehat{\Phi}) = \frac{1}{n} \sum_{i=1}^n \langle \Phi(\mathbf{x}^{(i)}), \widehat{\Phi} \rangle \Phi(\mathbf{x}^{(i)})$$

Therefore, equation (7.2) is equivalent to the following system of equations in $\alpha \in \mathbb{R}^n$:

$$\widehat{\Sigma}_n\left(\sum_{i=1}^n \alpha_i \Phi(\mathbf{x}^{(i)})\right) = \lambda \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}^{(i)})$$

For the above set of equations, we have

$$\text{LHS} = \frac{1}{n} \sum_{j=1}^n \sum_{i=1}^n \alpha_i \langle \Phi(\mathbf{x}^{(j)}), \Phi(\mathbf{x}^{(i)}) \rangle_{\mathcal{F}} \Phi(\mathbf{x}^{(j)})$$

Using the fact that $\langle \Phi(\mathbf{x}^{(j)}), \Phi(\mathbf{x}^{(i)}) \rangle_{\mathcal{F}} = \langle \Phi(\mathbf{x}^{(i)}), \Phi(\mathbf{x}^{(j)}) \rangle_{\mathcal{F}} = \mathbb{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$, the above system of equations may be written as

$$\frac{1}{n} \sum_{i,j=1}^n \alpha_i \mathbb{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \Phi(\mathbf{x}^{(j)}) = \lambda \sum_{i=1}^n \alpha_i \Phi(\mathbf{x}^{(i)})$$

Taking inner products with $\Phi(\mathbf{x}^{(l)}), l = 1, \dots, n$, we get

$$\frac{1}{n} \sum_{i,j=1}^n \alpha_i \mathbb{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \mathbb{K}(\mathbf{x}^{(j)}, \mathbf{x}^{(l)}) = \lambda \sum_{i=1}^n \alpha_i \mathbb{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(l)}).$$

We now have a set of n linear equations in the vector $\alpha \in \mathbb{R}^n$. In matrix-vector form, it can be written very simply as

$$K^2 \alpha = \lambda n K \alpha,$$

where $K \in \mathbb{R}^{n \times n}$ is the familiar kernel Gram matrix, with entries $K_{ij} = \mathbb{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$.

The only solutions of this equation that are of interest to us are those that satisfy

$$K \alpha = \lambda n \alpha.$$

This is simply an eigenvalue/eigenvector problem in the matrix K . □