# Kernel PCA : Mercer's Thm

Given $x_1, \ldots, x_m \in \mathbb{R}^n$ i.i.d. random vectors

$$\mu = \mathbb{E}[x] \quad \text{and} \quad C = \mathbb{E}\left[(x-\mu)(x-\mu)^T\right]$$

Principal Components of $x \in \mathbb{R}^N$ are entries of

$$y = U^T(x-\mu) = \begin{bmatrix} - & u_1^T & - \\ - & u_2^T & - \\ & \vdots & \\ - & u_N^T & - \end{bmatrix} (x-\mu)$$

**View 1**    $\mathbb{E}[y] = 0$ and $\mathbb{E}[yy^T] = \Lambda$

**View 2**    Direction $u_j$ maximizes $\text{var}(y^{(j)})$ s.t.
$$u_j \perp \{u_1, u_2, \ldots, u_{j-1}\}$$

$$u_j = \underset{\substack{\|v_j\|=1 \\ v_j \perp \{u_1, \ldots, u_{j-1}\}}}{\arg\max} \mathbb{E}\left[|v_j^T(x-\mu)|^2\right]$$

PCA can only "find" detect linear features in the data.

# Nonlinear Component Analysis

To identify "nonlinear" structure in the data, we can "add" new variables:

$$X_1 = \begin{pmatrix} X_1^{(1)} \\ X_1^{(2)} \end{pmatrix}, \quad \cdots, \quad X_m = \begin{pmatrix} X_m^{(1)} \\ X_m^{(2)} \end{pmatrix}$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

Add $\quad X_k^{(3)} = (X_k^{(1)})^2, \quad X_k^{(4)} = X_k^{(1)} X_k^{(2)}, \quad X_k^{(5)} = (X_k^{(2)})^2$

$$\Rightarrow \quad \text{Can find} \quad r^2 = (x^{(1)})^2 + (x^{(2)})^2$$

Let $\varphi = \mathbb{R}^n \to \mathbb{R}^d$ be a dictionary of features that "lift" the data into a higher-dimensional space ($d \gg n$).

$$\varphi(\overset{\in \mathbb{R}^n}{x}) = [\varphi_1(x), \varphi_2(x), \cdots, \varphi_d(x)]^T$$

$$\mu = \sum_{j=1}^{m} \varphi(x_j), \quad C = \frac{1}{m-1} \sum_{j=1}^{m} (\varphi(x_j) - \mu)(\varphi(x_j) - \mu)^T$$

We can run PCA in the new higher-dim feature space to detect nonlinear structure:

$$C = U \Lambda U^T \implies \Psi(x_j) = U^T \phi(x_j)$$

Diagonalize
$d \times d$ cov. matrix

principle components
of mapped data
in feature space

# Kernel PCA

$$C = \frac{1}{m-1} \sum_{i=1}^{m} (\phi(x_j) - \mu)(\phi(x_j) - \mu)^T = \frac{1}{m-1} \overset{d \times m}{B} \overset{m \times d}{B^T}$$

$$[\quad][\quad\quad]$$

To compute nonzero eigenpairs:

$$\lambda \neq 0 \quad \overset{d \times d}{Cu = \lambda u} \quad \underset{u = \frac{1}{\sqrt{m-1}} Bv}{\Longleftrightarrow} \quad \frac{1}{m-1} \overset{m \times m}{B^T B v = \lambda v}$$

To compute the principal components:

$$u^T(\phi(x_j) - \mu) = \frac{1}{\sqrt{m-1}} v^T B^T (\phi(x_j) - \mu)$$

$$= \frac{1}{\sqrt{m-1}} v^T B^T B$$

# The Kernel Matrix

To avoid working in the $d$-dimensional space

$$(B^T B)_{ij} = [\varphi(x_i) - \mu]^T [\varphi(x_j) - \mu]$$

$$= \sum_{k=1}^{d} (\varphi_k(x_j) - \mu)(\varphi_k(x_i) - \mu)$$

For simplicity, take $\mu = 0$ (not necessary):

$$(B^T B)_{ij} = \sum_{k=1}^{d} \varphi_k(x_j) \varphi_k(x_i)$$

Define the kernel $K: \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$

$$K(x,y) = \sum_{k=1}^{d} \varphi_k(x) \varphi_k(y)$$

Instead of starting w/ $\varphi: \mathbb{R}^n \to \mathbb{R}^d$, we start with a self-adjoint semi-definite kernel $K: \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$, continuous on Hilbert-Schmidt $\mathbb{R}^n \times \mathbb{R}^n$. Then, we have Mercer's Thm:

$$k(x,y) = \sum_{k=1}^{d \to \mathbb{N} \cup \{\infty\}} \lambda_n u_n(x) u_n(y)$$

Converges pointwise, absolutely & uniformly.

Mercer's theorem allows us to write

$$\int d^{\to \mathbb{N} \cup \{\infty\}}$$

$$k(x_i, x_j) = \sum_{k=1}^{d} \lambda_n u_n(x_i) u_n(x_j) = (B^T B)_{ji}$$

Dictionary $\quad \varphi(x) = \left[ \sqrt{\lambda_1} u_1(x), \sqrt{\lambda_2} u_2(x), \ldots, \sqrt{\lambda_d} u_d(x) \right]$

$$\varphi_1(x) \quad \varphi_2(x) \ldots \quad \varphi_d(x)$$

## Kernel PCA (mean $\mu = 0$)

Form $\quad (K)_{ij} = k(x_i, x_j) \qquad 1 \leq i, j \leq m$

Compute $\quad \frac{1}{\sqrt{M-1}} K v_\ell = \lambda v_\ell \qquad \ell = 1, 2, 3, \ldots$

Project $\quad c_\ell = \frac{1}{\sqrt{M-1}} v_\ell^T K \qquad \ell = 1, 2, 3, \ldots$
($\ell$th
(principal
component))

$\Rightarrow$ Implicitly run PCA in $\infty$-dim.
feature space via "kernel trick."

$\Rightarrow$ Kernel PCA "finds" nonlinear
structures in $\mathbb{R}^N$ that are well
approximated by kernel eigfuns

associated w/ largest eigenvalues

=> Different Kernels lead to different biases based on eigenvalue decay and "dominant" eigenspaces.