

Kernel Principal Component Analysis

Alex Beeny abeeny@siue.edu

Draft: March 14, 2024

Contents

1	Introduction	1
2	Principal Component Analysis	2
3	Reproducing Kernel Hilbert Space	6
4	Kernel PCA	12
5	Conclusion	13
A	Linear Algebra	13
B	Riesz Representation Theorem	15
C	Mercer's Theorem	15
D	Code	15

Abstract

Principal component analysis (PCA) and kernel methods are tools often used in data science. The underlying theory of these tools depend on the properties of a special type of Hilbert space called a reproducing kernel Hilbert space (RKHS). This paper explores the essence of RKHSs using data science examples, in particular, PCA and kernel PCA. When kernel methods are applied to PCA, we can analyze nonlinear data in a high-dimensional feature space with some nice properties.

1 Introduction

In linear regression, the equation of a line $y = a_0 + a_1x$ is used to model observations based on training data. Here, the input variable x is used to predict the response variable y . The parameters a_0 and a_1 are chosen such that

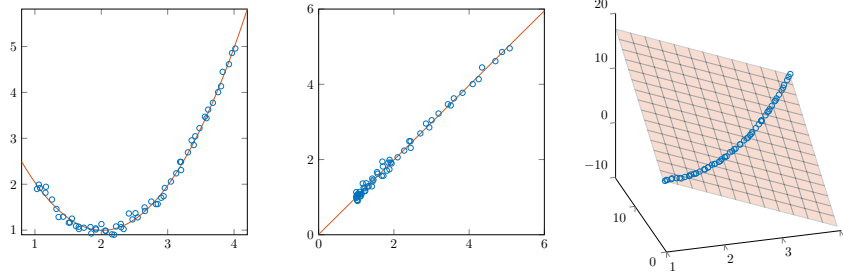


Figure 1: Plotting y against x (left) and the derived feature $z = \phi(x)$ (middle). The 3D plot (right) graphs the output y against x and x^2

the residual error¹ is minimized. It seems natural to model data using polynomial equations in a similar way, that is, determine a_0, a_1, \dots, a_n such that

$$y = a_0 + a_1x + a_2x^2 + \dots + a_nx^n \quad (1)$$

minimizes the residual error.

In multiple linear regression, the equation of a hyperplane

$$y = a_0 + a_1x_1 + a_2x_2 + \dots + a_nx_n \quad (2)$$

is used to predict y using inputs x_1, x_2, \dots, x_n . It follows that these observations are points in $n + 1$ dimensions.

Example 1.1. Suppose we have a set of points $\{(x^{(i)}, y^{(i)})\}_{i=1}^k$ in \mathbb{R}^2 . Using polynomial regression, we fit the model $y \sim 1 + x + x^2$. We can derive a new variable from x to get

$$z = \phi(x) = a_0 + a_1x + a_2x^2.$$

By transforming our original data, the quadratic relationship between x and y can be viewed as a linear relationship between z and y . This shows that polynomial regression is just a special case of multiple regression where each power of x is treated as a separate dimension. See Figure 1

Here, we make the distinction between different kinds of input variables. We say that x is an **attribute** and that z is a **feature**.

2 Principal Component Analysis

Principal component analysis (PCA) is a coordinate transform that minimizes projection residuals along each of the major axes e_1, e_2, \dots, e_m of the transformed space. These axes are referred to as **principal components**. After the

¹The residual for a given observation $(x^{(i)}, y^{(i)})$ is $|y^{(i)} - (a_0 + a_1x^{(i)})|$.

PCA transform is applied to a set of observations, the smallest projection error is along the first principal component. If we project down to two dimensions, the first two principal components have the smallest projection error. In general, projecting down to $k \leq m$ dimensions is most accurate along e_1, e_2, \dots, e_k .

Conversely, projecting along e_2, \dots, e_m will result in the largest residual error in $m - 1$ dimensions. It follows that e_1 must be the direction with the highest variance. Then e_2 is the direction with the next highest variance, and so on.

Example 2.1. Consider the following matrix

$$A = \begin{bmatrix} 5 & 3 & 6 & 7 & 6 \\ 4 & 5 & 7 & 1 & 3 \\ 5 & 7 & 6 & 1 & 0 \\ 6 & 10 & 12 & 12 & 11 \\ 9 & 10 & 12 & 13 & 9 \end{bmatrix}.$$

The column means are $\mu = [5.8, 7, 8.6, 6.8, 5.8]$. Then the mean-centered data becomes

$$X = A - \mu = \frac{1}{5} \begin{bmatrix} -4 & -20 & -13 & 1 & 1 \\ -9 & -10 & -8 & -29 & -14 \\ -4 & 0 & -13 & -29 & -29 \\ 1 & 15 & 17 & 26 & 26 \\ 16 & 15 & 17 & 31 & 16 \end{bmatrix}.$$

The covariance matrix is

$$C = X^T X = \frac{1}{5} \begin{bmatrix} 74 & 85 & 93 & 179 & 104 \\ 85 & 190 & 170 & 225 & 150 \\ 93 & 170 & 196 & 313 & 238 \\ 179 & 225 & 313 & 664 & 484 \\ 104 & 150 & 238 & 484 & 394 \end{bmatrix}.$$

Diagonalizing C gives

$$V = \begin{bmatrix} 0.1888 & -0.2020 & -0.6366 & 0.5495 & -0.4651 \\ 0.2755 & -0.7886 & 0.1472 & -0.4502 & -0.2791 \\ 0.3606 & -0.3464 & 0.3128 & 0.5836 & 0.5582 \\ 0.6979 & 0.2522 & -0.4422 & -0.3707 & 0.3411 \\ 0.5209 & 0.3922 & 0.5288 & 0.1316 & -0.5271 \end{bmatrix},$$

$$D = \begin{bmatrix} 264.8458 & 0 & 0 & 0 & 0 \\ 0 & 27.9766 & 0 & 0 & 0 \\ 0 & 0 & 9.3198 & 0 & 0 \\ 0 & 0 & 0 & 1.4579 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}.$$

If we keep all 5 principal component vectors, then $V_5 = V$ and the projection

of X along V is

$$P = XV = \begin{bmatrix} -1.9469 & 4.3453 & -0.8756 & -0.2039 & 0 \\ -6.9742 & -0.0660 & 1.4352 & 0.7590 & 0 \\ -8.1577 & -2.6752 & -0.8063 & -0.5704 & 0 \\ 8.4282 & -0.2330 & 1.8282 & -0.4996 & 0 \\ 8.6507 & -1.3711 & -1.5815 & 0.5149 & 0 \end{bmatrix}.$$

Here, the last column of P is the zero vector because the last eigenvalue of C is zero². To perfectly reconstruct A , we need $k = 4$ principal components and the row vector μ

$$A = PV^T + \mu = PV_4^T + \mu.$$

If we use $k = 3$ principal components, then the projection of X onto V_3 is

$$P = XV_3 = \begin{bmatrix} -1.9469 & 4.3453 & -0.8756 \\ -6.9742 & -0.0660 & 1.4352 \\ -8.1577 & -2.6752 & -0.8063 \\ 8.4282 & -0.2330 & 1.8282 \\ 8.6507 & -1.3711 & -1.5815 \end{bmatrix}$$

and A is approximately reconstructed by

$$A \approx PV_3^T + \mu = \begin{bmatrix} 5.1 & 2.9 & 6.1 & 6.9 & 6.0 \\ 3.6 & 5.3 & 6.6 & 1.3 & 2.9 \\ 5.3 & 6.7 & 6.3 & 0.8 & 0.1 \\ 6.3 & 9.8 & 12.3 & 11.8 & 11.1 \\ 8.7 & 10.2 & 11.7 & 13.2 & 8.9 \end{bmatrix}.$$

We can compute the reconstruction error using

$$E_k = \|A - (PV_k^T + \mu)\|_F,$$

where $\|\cdot\|_F$ is the Frobenius norm. By the SVD, we have $X = USV^T$, where $S = \sqrt{D}$. So, the projection of X onto V_k is

$$P = XV_k = US_k,$$

where S_k is the diagonal matrix of the first k singular values. Then the reconstruction error becomes

$$\begin{aligned} \|A - (PV_k^T + \mu)\|_F &= \|(A - \mu) - PV_k^T\|_F \\ &= \|X - PV_k^T\|_F \\ &= \|USV^T - US_kV^T\|_F \\ &= \|U(S - S_k)V^T\|_F \\ &= \|S - S_k\|_F \\ &= \sigma_k + \sigma_{k+1} + \cdots + \sigma_p. \end{aligned}$$

Hence,

$$E_3 = \sigma_3 + \sigma_4 = \sqrt{1.4579} + 0 = 1.2074.$$

²Since we subtracted the column means from a square matrix A , the dimension of the row space was reduced to 4.

2.1 Linear regression and PCA

[14] Let $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$ be observation vectors and $\mathbf{y} \in \mathbb{R}^n$ be a target vector. A linear regression model finds a vector of weights $\mathbf{w} = (w_1, w_2, \dots, w_n)$ to determine the linear function

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} = \sum_{i=1}^n w_i x_i, \quad (3)$$

where $\mathbf{x} = (x_1, x_2, \dots, x_n)$ is an input vector. The residual error of each observation is $\mathbf{y} - f(\mathbf{x}_i)$, for $i = 1, 2, \dots, m$. Then the residual sum of squares is given to be

$$RSS = \sum_{i=1}^m (\mathbf{y} - f(\mathbf{x}_i))^2 = \sum_{i=1}^m (\mathbf{y} - \mathbf{w} \cdot \mathbf{x}_i)^2 = (\mathbf{y} - X^\top \mathbf{w})^\top (\mathbf{y} - X^\top \mathbf{w}), \quad (4)$$

where $X = [\mathbf{x}_1, \dots, \mathbf{x}_m]$. By minimizing RSS , the magnitude of the residuals will be as small as possible, producing the optimal linear model $f(\mathbf{x})$. Therefore, this method is known as least squares regression. Differentiate (4) with respect to \mathbf{w} and set equal to zero so that

$$2X^\top \mathbf{y} - 2X^\top X \mathbf{w} = 0. \quad (5)$$

This leads to the normal equation $X^\top X \mathbf{w} = X^\top \mathbf{y}$. Thus,

$$\mathbf{w} = (X^\top X)^{-1} X^\top \mathbf{y} \quad (6)$$

gives the optimal linear model which minimizes residual error.

Example 2.2. In two dimensions, let $\mathbf{x} = (x_1, \dots, x_m)$ and $\mathbf{y} = (y_1, \dots, y_m)$. Then the least squares regression model is given by

$$f(x) = \frac{\text{cov}(\mathbf{x}, \mathbf{y})}{\text{var}(\mathbf{x})} (x - \bar{\mathbf{x}}) + \bar{\mathbf{y}}, \quad (7)$$

Since \mathbf{x} is treated as an input variable and \mathbf{y} is treated as an output, this is a type of supervised learning. In contrast, PCA is an unsupervised learning technique. PCA organizes the variables in the input space to reveal any patterns in the underlying data. If we have input variables \mathbf{x}_1 and \mathbf{x}_2 , we can use PCA to find a linear pattern among the inputs without trying to predict an output. First, we find the direction of the largest variance λ_{\max} using the covariance matrix

$$\begin{bmatrix} \text{cov}(\mathbf{x}_1, \mathbf{x}_1) & \text{cov}(\mathbf{x}_1, \mathbf{x}_2) \\ \text{cov}(\mathbf{x}_2, \mathbf{x}_1) & \text{cov}(\mathbf{x}_2, \mathbf{x}_2) \end{bmatrix} = \begin{bmatrix} \text{var}(\mathbf{x}_1) & \text{cov}(\mathbf{x}_1, \mathbf{x}_2) \\ \text{cov}(\mathbf{x}_1, \mathbf{x}_2) & \text{var}(\mathbf{x}_2) \end{bmatrix}.$$

In particular,

$$\lambda_{\max} = \frac{1}{2} \left(\text{var}(\mathbf{x}) + \text{var}(\mathbf{y}) + \sqrt{(\text{var}(\mathbf{x}) - \text{var}(\mathbf{y}))^2 + 4 \text{cov}(\mathbf{x}, \mathbf{y})^2} \right).$$

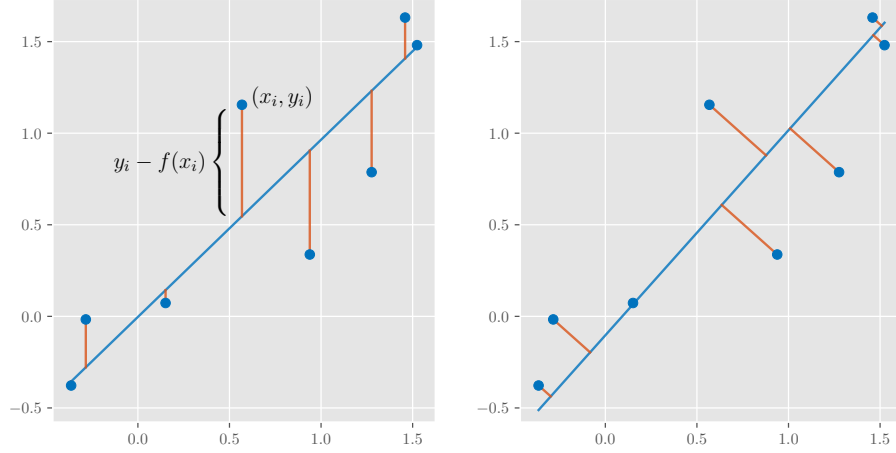


Figure 2: Least squares regression model $f(x)$ (left) minimizes the sum of squared errors while total least squares, i.e., the PCA model $g(x)$ (right) minimizes the orthogonal projections.

Then the PCA regression model becomes

$$g(x) = \frac{\lambda_{\max} - \text{var}(\mathbf{x}_1)}{\text{cov}(\mathbf{x}_1, \mathbf{x}_2)}(x - \bar{\mathbf{x}}_1) + \bar{\mathbf{x}}_2. \quad (8)$$

In this case, the projection residuals are orthogonal to the PCA regression line. See Figure 2.

3 Reproducing Kernel Hilbert Space

In this section, our goal is to establish properties of Hilbert spaces and kernel functions that can be used to modify the PCA algorithm. To begin, we will briefly mention definitions and results from analysis [7], [10] and matrix theory [5].

Definition 3.1. [15] Let X be a (real) vector space. An *inner product* is a function $\langle \cdot, \cdot \rangle : X \times X \rightarrow \mathbb{R}$ which satisfies the following properties:

1. Symmetry. For all $x, y \in X$,

$$\langle x, y \rangle = \langle y, x \rangle.$$

2. Linear in the first argument. For all $x, y, z \in X$, $\alpha, \beta \in \mathbb{R}$,

$$\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle.$$

3. Positive definite. For all $x \in X$,

$$\langle x, x \rangle \geq 0$$

and $\langle x, x \rangle = 0$ if and only if $x = 0$.

An *inner product space* is a vector space along with an inner product.

Since real inner products are symmetric and linear in the first argument,

$$\langle z, \alpha x + \beta y \rangle = \langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle = \alpha \langle z, x \rangle + \beta \langle z, y \rangle.$$

So, we also have linearity in the second argument.

The norm induced by an inner product is defined as

$$\|x\| = \langle x, x \rangle^{1/2}$$

and the metric induced by this norm is

$$d(x, y) = \|x - y\| = \langle x - y, x - y \rangle^{1/2}.$$

It follows that an inner product space is also a normed space and a metric space. So, the induced norm will have the following properties for all $x, y \in X$ and $\alpha \in \mathbb{R}$:

1. Triangle inequality. $\|x + y\| \leq \|x\| + \|y\|$;
2. Scalar multiplication. $\|\alpha x\| = |\alpha| \|x\|$;
3. Positivity. $\|x\| \geq 0$ and $\|x\| = 0$ if and only if $x = 0$.

Definition 3.2 (Hilbert space). A Hilbert space is a complete inner product space. For a Hilbert space H , we sometimes denote the inner product as $\langle \cdot, \cdot \rangle_H$ to avoid ambiguity.

Two Hilbert spaces H and L (over the same field) are said to be *isomorphic* if there is a bijection $T : H \rightarrow L$ such that

$$\langle x, y \rangle_H = \langle Tx, Ty \rangle_L, \tag{9}$$

for every $x, y \in H$. In [7], Kreyszig shows that two Hilbert spaces are isomorphic if and only if they have the same dimension. If V is an inner product space that is not complete, then it can be extended to a Hilbert space by completion. The completion of an inner product space is denoted as \overline{V} and is unique up to isomorphism.

A Hilbert space is said to be *separable* if it contains a dense countable subset. It can be shown [7] that a Hilbert space is separable if and only if it has a countable orthonormal basis. The following example demonstrates a useful property of separable Hilbert spaces.

Example 3.3. [10] The space of square-summable (real) sequences is defined as

$$\ell^2(A) = \left\{ x : A \rightarrow \mathbb{R} \mid \sum_{a \in A} x_a^2 < \infty \right\}. \quad (10)$$

Given the inner product

$$\langle x, y \rangle = \sum_{a \in A} x_a y_a, \quad (11)$$

$\ell^2(A)$ is a Hilbert space. Moreover, $\ell^2(A)$ is separable if and only if A is countable. It follows that the sequence space $\ell^2 = \ell^2(\mathbb{N})$ is the separable Hilbert space of square-summable sequences. Due to the Riesz-Fischer theorem, every infinite-dimensional Hilbert space is isomorphic to ℓ^2 .

Definition 3.4 (Gram matrix). [5] Let x_1, x_2, \dots, x_n be vectors in an inner product space X . The matrix $G \in \mathbb{R}^{n \times n}$ is called a *Gram matrix* if $[G]_{ij} = \langle x_i, x_j \rangle$.

Example 3.5. Consider the vectors in \mathbb{R}^3 :

$$\mathbf{v}_1 = \begin{bmatrix} v_{11} \\ v_{21} \\ v_{31} \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} v_{12} \\ v_{22} \\ v_{32} \end{bmatrix}, \quad \mathbf{v}_3 = \begin{bmatrix} v_{13} \\ v_{23} \\ v_{33} \end{bmatrix}, \quad \mathbf{v}_4 = \begin{bmatrix} v_{14} \\ v_{24} \\ v_{34} \end{bmatrix}.$$

The Gram matrix for these vectors is

$$G = \begin{bmatrix} \mathbf{v}_1^\top \mathbf{v}_1 & \mathbf{v}_1^\top \mathbf{v}_2 & \mathbf{v}_1^\top \mathbf{v}_3 & \mathbf{v}_1^\top \mathbf{v}_4 \\ \mathbf{v}_2^\top \mathbf{v}_1 & \mathbf{v}_2^\top \mathbf{v}_2 & \mathbf{v}_2^\top \mathbf{v}_3 & \mathbf{v}_2^\top \mathbf{v}_4 \\ \mathbf{v}_3^\top \mathbf{v}_1 & \mathbf{v}_3^\top \mathbf{v}_2 & \mathbf{v}_3^\top \mathbf{v}_3 & \mathbf{v}_3^\top \mathbf{v}_4 \\ \mathbf{v}_4^\top \mathbf{v}_1 & \mathbf{v}_4^\top \mathbf{v}_2 & \mathbf{v}_4^\top \mathbf{v}_3 & \mathbf{v}_4^\top \mathbf{v}_4 \end{bmatrix}.$$

If V is a matrix whose columns are $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \mathbf{v}_4$, then we can write $G = V^\top V$.

Definition 3.6 (kernel). Let X be a nonempty set and $k : X \times X \rightarrow \mathbb{R}$. Then k is a (positive semi-definite) *kernel* if

1. k is symmetric: for all $x, y \in X$,

$$k(x, y) = k(y, x);$$

2. k is positive semi-definite: for all $n \in \mathbb{N}$, if $x_1, \dots, x_n \in X$ and $c_1, \dots, c_n \in \mathbb{R}$, then

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \geq 0.$$

Equivalently, the matrix $K \in \mathbb{R}^{n \times n}$ whose entries are $K_{ij} = k(x_i, x_j)$ is positive semi-definite, that is, $w^\top K w \geq 0$ for all $w \in \mathbb{R}^n$.

Definition 3.7 (feature map). [9] Let F be a Hilbert space of functions $f : X \rightarrow \mathbb{R}$. A feature map is a function $\Phi : X \rightarrow F$.

Our goal is to have kernels written as inner products of feature maps:

$$k(x, y) = \langle \Phi(x), \Phi(y) \rangle_H.$$

If Φ is a known feature map, then k would surely be symmetric and positive definite since inner products are. If instead we know k is a kernel, then we need to show that there is a unique Hilbert space with the desired inner product.

To do:

- move kernel/Gram matrix to its own definition
- Define linear functionals
- Define dual space (maybe?)
- Riesz representation theorem
- Define evaluation functionals
- Define RKHS (continuous evaluation functionals)
- An RKHS defines a unique reproducing kernel (by Riesz representation theorem).
- Mercer's theorem.
- A kernel defines a feature map.
- A kernel defines a unique RKHS (by Mercer/Moore-Aronszajn).
- We finally get $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_H$.

Theorem 3.8. [9, 14] Suppose k_1 and k_2 are kernels over $X \times X$. The following functions kernels.

1. $k(x, y) = a_1 k_1(x, y) + a_2 k_2(x, y)$ for all $a_1, a_2 \geq 0$.
2. $k(x, y) = k_1(x, y) k_2(x, y)$.
3. $k(x, y) = a_0 + a_1 k_1(x, y) + a_2 k_1(x, y)^2 + \cdots + a_n k_1(x, y)^n$ for all $n \in \mathbb{N}$ and $a_0, \dots, a_n \geq 0$.
4. $k(x, y) = k_1(h(x), h(y))$ for all $h : X \rightarrow X$.
5. $k(x, y) = g(x)g(y)$ for all $g : X \rightarrow \mathbb{R}$.
6. $k(x, y) = \exp(k_1(x, y))$.

Proof. Let $x_1, \dots, x_n \in X$ and $c_1, \dots, c_n \in \mathbb{R}$.

1. Let $k = a_1 k_1 + a_2 k_2$ for $a_1, a_2 \geq 0$. Since k_1 and k_2 are symmetric,

$$k(x, y) = a_1 k_1(x, y) + a_2 k_2(x, y) = a_1 k_1(y, x) + a_2 k_2(y, x) = k(y, x),$$

for all $x, y \in X$. So, k is symmetric.

Since k_1 and k_2 are positive semi-definite and $a_1, a_2 \geq 0$,

$$\begin{aligned} \sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j (a_1 k_1(x_i, x_j) + a_2 k_2(x_i, x_j)) \\ &= a_1 \sum_{i=1}^n \sum_{j=1}^n c_i c_j k_1(x_i, x_j) + a_2 \sum_{i=1}^n \sum_{j=1}^n c_i c_j k_2(x_i, x_j) \\ &\geq 0. \end{aligned}$$

So, k is positive semi-definite.

2. Let $k = k_1 k_2$. Define K so that $[K]_{ij} = k(x_i, x_j) = k_1(x_i, x_j) k_2(x_i, x_j)$. Let K_1 and K_2 be the Gram matrices for k_1 and k_2 , respectively. Then K_1, K_2 have orthonormal eigenvectors and nonnegative eigenvalues such that

$$\begin{aligned} K_1 &= V L V^\top \\ &= \begin{bmatrix} v_{11} & \cdots & v_{1n} \\ \vdots & \ddots & \vdots \\ v_{n1} & \cdots & v_{nn} \end{bmatrix} \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{bmatrix} \begin{bmatrix} v_{11} & \cdots & v_{n1} \\ \vdots & \ddots & \vdots \\ v_{1n} & \cdots & v_{nn} \end{bmatrix} \\ &= \begin{bmatrix} \sum_{j=1}^n \lambda_j v_{1j} v_{1j} & \cdots & \sum_{j=1}^n \lambda_j v_{nj} v_{1j} \\ \vdots & \ddots & \vdots \\ \sum_{j=1}^n \lambda_j v_{1j} v_{nj} & \cdots & \sum_{j=1}^n \lambda_j v_{nj} v_{nj} \end{bmatrix} \\ &= \sum_{j=1}^n \lambda_j \begin{bmatrix} v_{1j} v_{1j} & \cdots & v_{nj} v_{1j} \\ \vdots & \ddots & \vdots \\ v_{1j} v_{nj} & \cdots & v_{nj} v_{nj} \end{bmatrix} \end{aligned}$$

and

$$K_2 = U M U^\top = \sum_{j=1}^n \mu_j \begin{bmatrix} u_{1j} u_{1j} & \cdots & u_{nj} u_{1j} \\ \vdots & \ddots & \vdots \\ u_{1j} u_{nj} & \cdots & u_{nj} u_{nj} \end{bmatrix}.$$

Let $\mathbf{v}_i = [v_{1i} \ \cdots \ v_{ni}]^\top$ and $\mathbf{u}_j = [u_{1j} \ \cdots \ u_{nj}]$, for all $i, j =$

$1, 2, \dots, n$. Then

$$\begin{aligned}
K &= K_1 \circ K_2 \\
&= \sum_{i=1}^n \lambda_i \begin{bmatrix} v_{1i}v_{1i} & \cdots & v_{ni}v_{1i} \\ \vdots & \ddots & \vdots \\ v_{1i}v_{ni} & \cdots & v_{ni}v_{ni} \end{bmatrix} \circ \sum_{j=1}^n \mu_j \begin{bmatrix} u_{1j}u_{1j} & \cdots & u_{nj}u_{1j} \\ \vdots & \ddots & \vdots \\ u_{1j}u_{nj} & \cdots & u_{nj}u_{nj} \end{bmatrix} \\
&= \sum_{i=1}^n \sum_{j=1}^n \lambda_i \mu_j \begin{bmatrix} v_{1i}u_{1j}v_{1i}u_{1j} & \cdots & v_{1i}u_{1j}v_{ni}u_{nj} \\ \vdots & \ddots & \vdots \\ v_{ni}u_{nj}v_{1i}u_{1j} & \cdots & v_{ni}u_{nj}v_{ni}u_{nj} \end{bmatrix} \\
&= \sum_{i=1}^n \sum_{j=1}^n \lambda_i \mu_j \begin{bmatrix} v_{1i}u_{1j} \\ \vdots \\ v_{ni}u_{nj} \end{bmatrix} [v_{1i}u_{1j} \quad \cdots \quad v_{ni}u_{nj}] \\
&= \sum_{i=1}^n \sum_{j=1}^n \lambda_i \mu_j (\mathbf{v}_i \circ \mathbf{u}_j)(\mathbf{v}_i \circ \mathbf{u}_j)^\top,
\end{aligned}$$

where \circ is the Hadamard product. Each $(\mathbf{v}_i \circ \mathbf{u}_j)(\mathbf{v}_i \circ \mathbf{u}_j)^\top$ is a symmetric positive semi-definite matrix. Since K_1, K_2 are positive semi-definite, we have $\lambda_i, \mu_i > 0$. Then K is symmetric positive semi-definite.

3. By part 2, k_1, k_1^2, \dots, k_1^n are kernels. By part 1, $a_0 + a_1 k_1 + a_2 k_1^2 + \dots + a_n k_1^n$ is a kernel.

4. Since $y_i = h(x_i) \in X$ for all $i = 1, 2, \dots, n$, we have

$$\begin{aligned}
\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j k_1(h(x_i), h(x_j)) \\
&= \sum_{i=1}^n \sum_{j=1}^n c_i c_j k_1(y_i, y_j) \\
&\geq 0.
\end{aligned}$$

5. Let $g : X \rightarrow \mathbb{R}$ and let $c_i g(x_i) = y_i \in \mathbb{R}$. If $k(x, y) = g(x)g(y)$, then

$$\begin{aligned}
\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) &= \sum_{i=1}^n \sum_{j=1}^n c_i g(x_i) c_j g(x_j) \\
&= \sum_{i=1}^n \sum_{j=1}^n y_i y_j \\
&= \left(\sum_{i=1}^n y_i \right)^2 \\
&\geq 0.
\end{aligned}$$

6. Let K_1 be the Gram matrix for k_1 . If $K_1 v = \lambda v$, then $K_1^m = \lambda^m v$ for all $m \in \mathbb{N}$. So,

$$(\exp K_1)v = \sum_{m=0}^{\infty} \frac{K_1^m v}{m!} = \sum_{m=0}^{\infty} \frac{\lambda^m v}{m!} = e^\lambda v.$$

Then $K = \exp K_1$ has eigenvalues e^λ . Since K_1 is positive semi-definite, it has real eigenvalues so that $e^\lambda > 0$. It follows that K is positive definite. \square

Theorem 3.9 (Gaussian kernel). *The function $k : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ defined by*

$$k(x, y) = \exp\left(\frac{-\|x - y\|_2^2}{\sigma^2}\right),$$

is a kernel.

4 Kernel PCA

Recall that linear PCA finds new components that reveal more information about the structure of high-dimensional data. Since PCA is an orthogonal projection, the original data is rotated within the original space of input variables. The work of Schölkopf, Smola, and Müller [11, 12] generalized PCA based on the successful application of kernel methods in support vector machines. In kernel PCA, the inner product of the input space is replaced with the inner product of a feature space. As such, the principal components of kernel PCA are nonlinear transformations of input variables, or features.

Using results from the previous section, a kernel function k defines a unique reproducing kernel Hilbert space H_k and feature map $\Phi : X \rightarrow H_k$ such that

$$k(x, y) = \langle \Phi(x), \Phi(y) \rangle, \quad (12)$$

for all $x, y \in H_k$. By replacing the dot product in the formulation of the PCA algorithm,

4.1 Centering in the feature space

Let $\Phi : X \rightarrow H_k$ be a feature map determined by a kernel k . Since Φ may be nonlinear, the image $\Phi(x)$ of a centered vector $x \in X$ is not guaranteed to be centered. For an effective PCA algorithm, it is necessary to compute the Gram matrix of centered vectors in the feature space. [12]

Given $x_1, \dots, x_m \in X$, the points

$$\Phi_0(x_i) = \Phi(x_i) - \frac{1}{m} \sum_{i=1}^m \Phi(x_i), \quad \text{for } i = 1, \dots, m \quad (13)$$

are the centered feature vectors in H_k . Then the centered Gram matrix becomes

$$\begin{aligned}
[K_0]_{ij} &= \langle \Phi_0(x_i), \Phi_0(x_j) \rangle \\
&= \left\langle \Phi(x_i) - \frac{1}{m} \sum_{p=1}^m \Phi(x_p), \Phi(x_j) - \frac{1}{m} \sum_{q=1}^m \Phi(x_q) \right\rangle \\
&= \langle \Phi(x_i), \Phi(x_j) \rangle - \frac{1}{m} \sum_{p=1}^m \langle \Phi(x_p), \Phi(x_j) \rangle \\
&\quad - \frac{1}{m} \sum_{q=1}^m \langle \Phi(x_i), \Phi(x_q) \rangle \\
&\quad + \frac{1}{m^2} \sum_{p=1}^m \sum_{q=1}^m \langle \Phi(x_p), \Phi(x_q) \rangle \\
&= k(x_i, x_j) - \frac{1}{m} \sum_{p=1}^m k(x_p, x_j) - \frac{1}{m} \sum_{q=1}^m k(x_i, x_q) + \frac{1}{m^2} \sum_{p=1}^m \sum_{q=1}^m k(x_p, x_q).
\end{aligned}$$

If K is the uncentered Gram matrix, then the formula can be written as

$$K_0 = K - \mathbf{1}_m K - K \mathbf{1}_m + \mathbf{1}_m K \mathbf{1}_m, \quad (14)$$

where $\mathbf{1}_m$ is the $m \times m$ matrix whose entries are $1/m$.

Example 4.1. Consider the problem of classifying points based on their radii. These points cannot be separated using a linear classifier in the two dimensions. However, by mapping them to a three-dimensional space, they can be separated by planes. Applying kernel PCA, these points can be sent to the RKHS associated with a Gaussian kernel without using an explicit feature map. The points in this high-dimensional feature space can then be projected onto the first three principal components to find separation boundaries. See Figure 3.

5 Conclusion

A Linear Algebra

Lemma A.1. *Symmetric matrices have orthogonal eigenvectors and real eigenvalues.*

Proof. □

Lemma A.2. *Positive semi-definite matrices have nonnegative eigenvalues.*

Proof. □

Lemma A.3. *Let A be an $m \times n$ matrix. Then $A^\top A$ is symmetric positive semi-definite.*

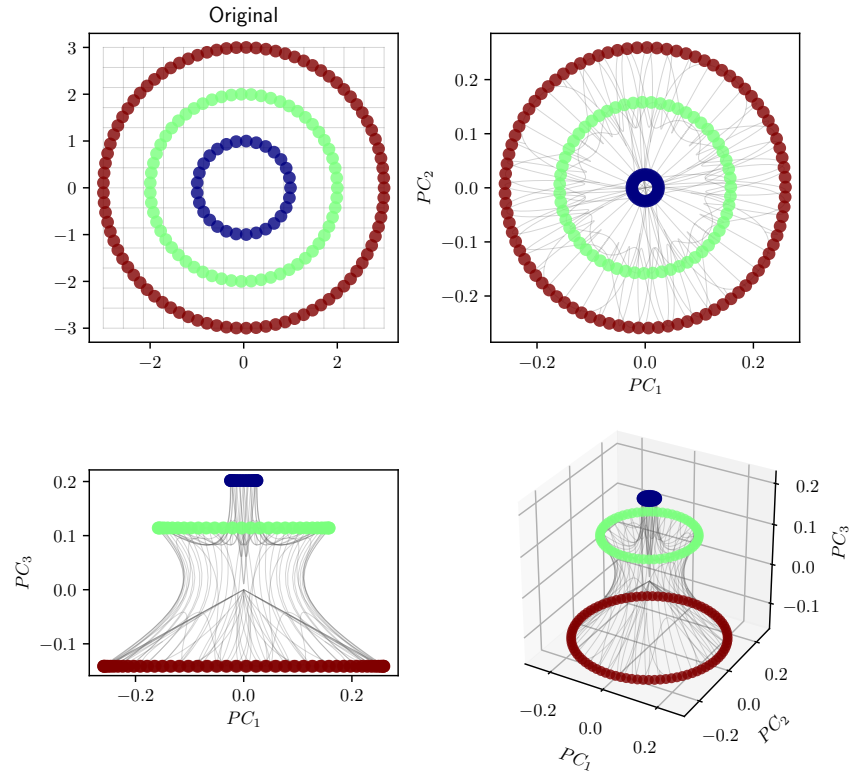


Figure 3: An idealized set of points in the plane are classified based on their radius. Kernel PCA with the Gaussian kernel is applied to find separation boundaries using the first three principal components.

Proof. First, $(A^\top A)^\top = A^\top (A^\top)^\top = A^\top A$. So $A^\top A$ is symmetric. \square

Definition A.4. Let A be an $m \times n$ matrix whose columns $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \in \mathbb{R}^m$ represent random variables and rows represent observations. Then the covariance matrix of A is given by

$$\text{cov}(A) = \begin{bmatrix} \text{cov}(\mathbf{a}_1, \mathbf{a}_1) & \text{cov}(\mathbf{a}_1, \mathbf{a}_2) & \cdots & \text{cov}(\mathbf{a}_1, \mathbf{a}_n) \\ \text{cov}(\mathbf{a}_2, \mathbf{a}_1) & \text{cov}(\mathbf{a}_2, \mathbf{a}_2) & \cdots & \text{cov}(\mathbf{a}_2, \mathbf{a}_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(\mathbf{a}_n, \mathbf{a}_1) & \text{cov}(\mathbf{a}_n, \mathbf{a}_2) & \cdots & \text{cov}(\mathbf{a}_n, \mathbf{a}_n) \end{bmatrix}.$$

If A is centered, i.e., the columns all have mean zero, then we can write

$$\text{cov}(A) = \frac{A^\top A}{m-1}.$$

B Riesz Representation Theorem

Theorem B.1 (Riesz Representation Theorem). [15] *Let $\phi : H \rightarrow \mathbb{R}$ be a continuous linear functional defined on a Hilbert space H . Then there exists a unique element $g \in H$ such that $\phi(g) = \langle f, g \rangle_H$ for all $g \in H$.*

Proof. \square

C Mercer's Theorem

D Code

References

- [1] Peter Bartlett. CS 281B / Stat 241B Statistical Learning Theory. Lecture notes, Spring 2008.
- [2] Peter Bartlett. CS 281B / Stat 241B Statistical Learning Theory. Lecture notes, Spring 2014.
- [3] Ji Chen and Xiaodong Li. Model-free nonconvex matrix completion: Local minima analysis and applications in memory-efficient kernel pca. *Journal of Machine Learning Research*, 20(142):1–39, 2019.
- [4] Thomas Hofmann, Bernhard Schölkopf, and Alexander J. Smola. Kernel methods in machine learning. *The Annals of Statistics*, 36(3), jun 2008.
- [5] R.A. Horn and C.R. Johnson. *Matrix Analysis*. Cambridge University Press, 2013.
- [6] The MathWorks Inc. Statistics and machine learning toolbox, 2022.

- [7] E. Kreyszig. *Introductory Functional Analysis with Applications*. Wiley Classics Library. Wiley, 1991.
- [8] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. Adaptive Computation and Machine Learning series. MIT Press, 2012.
- [9] Cynthia Rudin. *Intuition for the Algorithms of Machine Learning*. Self-pub, 2020.
- [10] Walter Rudin. *Real and Complex Analysis*. Higher Mathematics Series. McGraw-Hill Education, 1987.
- [11] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component analysis. In *International conference on artificial neural networks*, pages 583–588. Springer, 1997.
- [12] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998.
- [13] Cosma Rohilla Shalizi. *Advanced Data Analysis from an Elementary Point of View*. Draft textbook, 2021.
- [14] John Shawe-Taylor and Nello Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- [15] Christopher G. Small and D.L. Mcleish. *Hilbert Space Methods in Probability and Statistical Inference*. John Wiley & Sons, Inc, 1994.