

# Question on Kernel\_PCA: Sp24\_567\_VATSAL

sbmohant

April 2024

## 1 Introduction

OK

## 2 Kernel PCA

(12 points)

Just like the other linear algorithms discussed in class, it can often be useful to do PCA in a higher-dimensional feature space. In fact, we can actually apply the kernel trick to PCA! In this problem we will derive some properties of *Kernel PCA*.

Suppose we have  $n$  data points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ , and some feature map  $\phi(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}^M$ . Consider the data matrix  $\mathbf{X}$  of the original data, and the feature matrix  $\Phi$  of the transformed datapoints defined as usual as follows:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} \in \mathbb{R}^{n \times d}, \quad \Phi = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \\ \phi(\mathbf{x}_2)^\top \\ \vdots \\ \phi(\mathbf{x}_n)^\top \end{bmatrix} \in \mathbb{R}^{n \times M}.$$

Assume that  $\Phi$  is centered, so  $\mathbf{C} = \Phi^\top \Phi$  is the covariance matrix of the data in the high-dimensional feature space. Let  $\mathbf{K} = \Phi \Phi^\top$  be the Kernel matrix for the data with the feature transformation  $\phi(\cdot)$ . Recall that for any matrix  $\mathbf{A}$ ,  $(\lambda, \mathbf{v})$  is any eigenpair of  $\mathbf{A}$  if  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ . Since if  $\mathbf{v}$  is an eigenvector of  $\mathbf{A}$  then  $\alpha\mathbf{v}$  is also an eigenvector for any  $\alpha \neq 0$ , we usually normalize eigenvectors to have unit norm, i.e.  $\|\mathbf{v}\|_2 = 1$ .

(a) For any  $\lambda \neq 0$ , show that if  $(\lambda, \mathbf{v})$  is an eigenpair of  $\mathbf{C}$ , then  $(\lambda, \Phi\mathbf{v})$  is an eigenpair of  $\mathbf{K}$ . (2 points)

(b) Similarly for any  $\lambda \neq 0$ , show that if  $(\lambda, \mathbf{a})$  is an eigenpair of  $\mathbf{K}$ , then  $(\lambda, \Phi^\top \mathbf{a})$  is an eigenpair of  $\mathbf{C}$ . (2 points)

(c) By the previous two parts, for any  $\lambda \neq 0$ ,  $\lambda$  is an eigenvalue of  $\mathbf{C}$  if and only if  $\lambda$  is an eigenvalue of  $\mathbf{K}$ . This also means that the largest magnitude eigenvalue of  $\mathbf{K}$  is also the largest magnitude eigenvalue of  $\mathbf{C}$ .

Let  $\mathbf{a}_1$  be the eigenvector of  $\mathbf{K}$  corresponding to the largest magnitude eigenvalue  $\lambda_1$ . Show that the first principal component is given by,

$$\mathbf{v}_1 = \frac{\Phi^\top \mathbf{a}_1}{\|\Phi^\top \mathbf{a}_1\|_2}.$$

Assume that  $\|\mathbf{a}_1\|_2 = 1$ , and show that this expression can be simplified to,

$$\mathbf{v}_1 = \frac{\mathbf{\Phi}^\top \mathbf{a}_1}{\sqrt{\lambda_1}}. \quad (1)$$

(3 points)

(d) Given a datapoint  $\mathbf{x}$ , show that its projection  $\alpha_1$  onto the first principle component  $\mathbf{v}_1$  is given by

$$\alpha_1 = \frac{1}{\sqrt{\lambda_1}} \sum_{i=1}^n k(\mathbf{x}, \mathbf{x}_i) (\mathbf{a}_1)_i.$$

Here  $(\mathbf{a}_1)_i$  denotes the  $i$ -th coordinate of the vector  $\mathbf{a}_1$ . To derive this, use the expression for  $\mathbf{v}_1$  from Eq. 1 above. (Hint: The projection of a vector  $\mathbf{a}$  onto  $\mathbf{b}$  is  $\frac{\mathbf{b}^\top \mathbf{a}}{\|\mathbf{b}\|_2}$ .) (3 points)

(e) Figure 8a represents a dataset that doesn't seem to be amenable to PCA. So we apply a non-linear transformation  $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  given by

$$\phi \left( \begin{bmatrix} x \\ y \end{bmatrix} \right) = \begin{bmatrix} \sqrt{x^2 + y^2} \\ \arctan2(y, x) \end{bmatrix}$$

where the function  $\arctan2$  returns the angle from positive x-axis of the point (x,y) in the range  $(0, 2\pi)$ .

Similarly we define the inverse of  $\phi$  as  $\phi^{-1} \left( \begin{bmatrix} r \\ \theta \end{bmatrix} \right) = \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix}$

Locate the arrow(s) which best describe the following phenomena -

1) PCA does not always reduce dimensions but can find an alternative basis to represent the data.

2) Projecting the data onto the leading principal components are not always better for downstream classification tasks. Sometimes projecting onto the less significant principal components which contain lesser variance might be more suitable for downstream classification.

3) Kernel PCA can find a suitable non-linear basis to represent the original data on a manifold instead of a linear subspace as achieved by linear PCA.

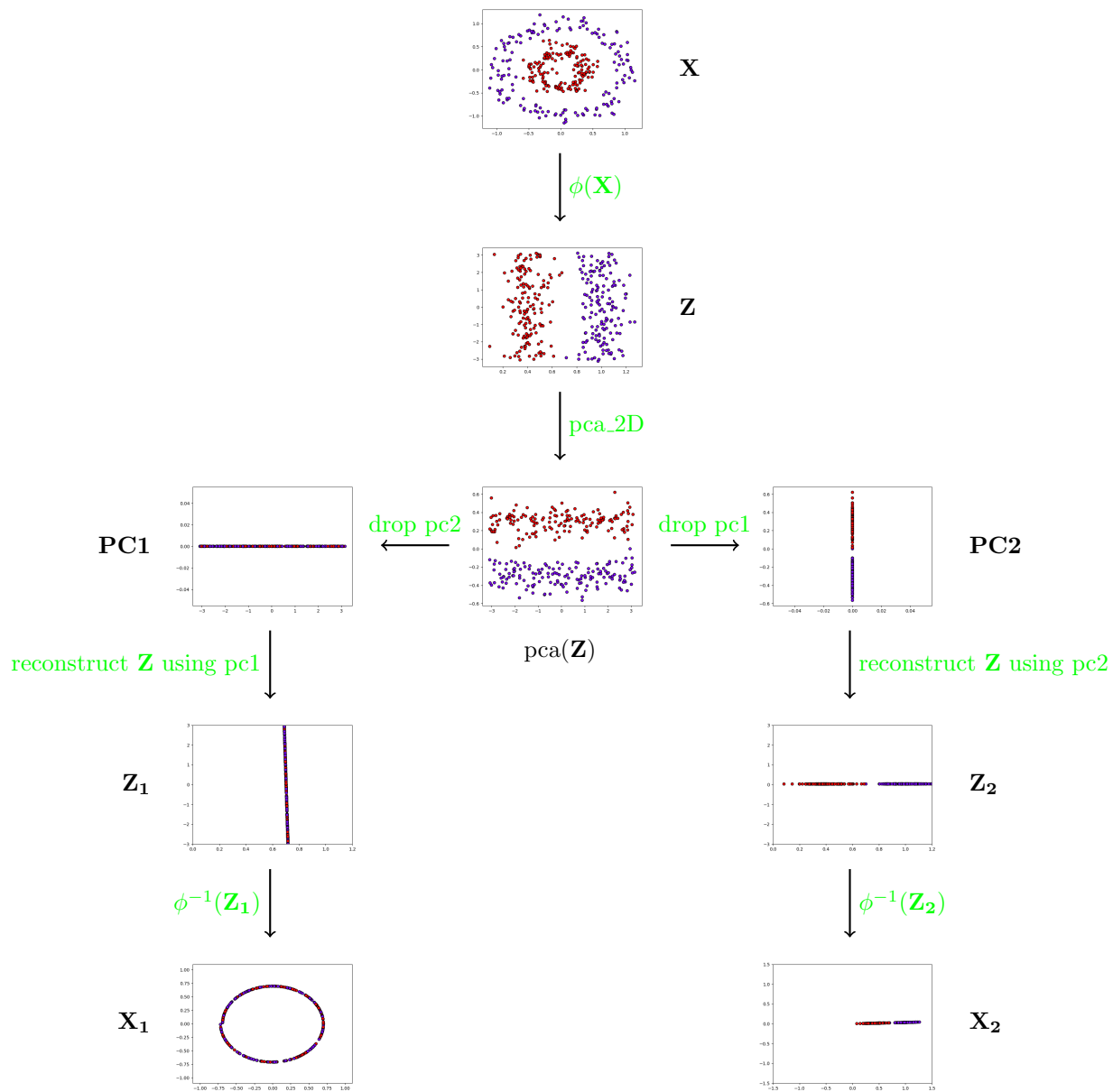


Figure 1: Kernel PCA step by step