

# Unsupervised Learning, HW II (The REMIX!)

## PCA and KPCA

Ofer Lipman, 201510435 and Daniel Shterenberg, 305199507

August 2021

## 1 PCA

### 1.1 Eigendecomposition

#### 1.1.1

Let  $A \in \mathbb{R}^{d \times d}$  where  $A = V\Lambda V^{-1}$ . We need to prove that  $\text{Tr}\{A\} = \sum_{i=1}^d \lambda_i(A)$ .

Since  $\Lambda$  is diagonal, then  $V\Lambda V^{-1}$  is an EVD and

$$\text{Tr}(A) = \text{Tr}(V\Lambda V^{-1}) = \text{Tr}(V^{-1}V\Lambda) = \text{Tr}(\Lambda)$$

Since  $\text{Tr}(A) = \text{Tr}(\Lambda)$  and  $\Lambda$  is a diagonal matrix (part of the EVD), then  $\Lambda[i, i]$  is the  $i$ 'th eigenvalue. ■

#### 1.1.2

We have  $A, B \in \mathbb{R}^{d \times d}$  such that  $A$  is diagonalizable and  $A \sim B$ .

Since  $A$  is diagonalizable, it has an EVD:  $U^t \Lambda U$  such that  $U$  consists of orthonormal vectors and  $\Lambda$  is a diagonal matrix.

In sections 1.1.1, we proved that  $\Lambda$  diagonal is the set of the matrix eigenvalues.

Now let's take  $B$ .

Since  $A \sim B$ , then exists an invertible matrix  $P$  such that  $B = PAP^{-1}$ .

We can see that  $B = PAP^{-1} = PU^t \Lambda UP^{-1}$ .

Denoting  $Z = U^t P$ , we will get  $B = Z\Lambda Z^{-1}$ , that means that  $Z\Lambda Z$  is EVD of  $B$ , and following section 1.1.1,  $\Lambda$  diagonal is the set of the eigenvalues.

Therefore  $\{\lambda_i(A)\}_{i=1}^d = \{\lambda_i(B)\}_{i=1}^d$ . ■

#### 1.1.3

Let  $A$  be a symmetric matrix ( $A \in \mathbb{R}^{d \times d}, A = A^T$ ). We would like to prove that

$$\lambda_i(A) > 0 \iff v^t A v > 0, \forall v \neq 0$$

$\implies \forall i : \lambda_i(A) > 0$ :

Since  $A$  is symmetric, it has an EVD,  $A = U\Sigma U^T$ .

All of  $A$ 's eigenvalues are larger than 0, then  $A$  has  $d$  eigenvectors  $u_1, \dots, u_d$  which represents a vectoric base for  $\mathbb{R}^d$ . Let  $v$  be some vector in  $\mathbb{R}^d$ . We can represent  $v$  as a linear combination of  $u_1, \dots, u_d$ :

$$v = \alpha_1 u_1 + \dots + \alpha_d u_d$$

Now, lets take a look at  $v^T Av$ :

$$\begin{aligned}
v^T Av &= (\alpha_1 \cdot u_1 + \dots + \alpha_d \cdot u_d)^T \cdot U \Sigma U^T \cdot (\alpha_1 \cdot u_1 + \dots + \alpha_d \cdot u_d) \\
&= (\alpha_1 \cdot u_1^T + \dots + \alpha_d \cdot u_d^T) \cdot U \Sigma U^T \cdot (\alpha_1 \cdot u_1 + \dots + \alpha_d \cdot u_d) \\
&= \begin{pmatrix} \alpha_1 \|u_1\|^2 & \alpha_2 \|u_2\|^2 & \dots & \alpha_d \|u_d\|^2 \end{pmatrix} \cdot \Sigma \cdot \begin{pmatrix} \alpha_1 \|u_1\|^2 \\ \alpha_2 \|u_2\|^2 \\ \vdots \\ \alpha_d \|u_d\|^2 \end{pmatrix} \\
&= \alpha_1^2 \lambda_1 + \dots + \alpha_d^2 \lambda_d
\end{aligned}$$

And since  $\lambda_1, \dots, \lambda_d > 0$ , we get  $v^T Av = \alpha_1^2 \lambda_1 + \dots + \alpha_d^2 \lambda_d > 0$  ■

$\Leftarrow$ :  $\forall v \neq 0 \in \mathbb{R}^d, v^T Av > 0$ :

Lets consider A's eigenvector  $u_i$  which holds  $Au_i = \lambda_i u_i$ .  $u_i^T Au_i = u_i^T (\lambda_i u_i) = \lambda_i (u_i^T u_i)$  and since  $u_i^T u_i \geq 0$  and  $u_i^T Au_i > 0$  we will get that the eigenvalue  $\lambda_i > 0$ . Of course that holds for any eigenvector  $u_i$  and therefore each eigenvalue  $\lambda_i > 0$  ■

## 1.2 PCA

### Full PCA

- Consider the data  $\mathcal{X} = \{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$  with mean  $\boldsymbol{\mu}_x \in \mathbb{R}^D$  and covariance  $\boldsymbol{\Sigma}_x \in \mathbb{R}^{D \times D}$ .
- Let  $\boldsymbol{\Sigma}_x = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T$  be the eigendecomposition of  $\boldsymbol{\Sigma}_x$ .
- Let  $\mathbf{z}_i = \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x)$

#### 1.2.1

Prove that:

1. The mean of  $\mathcal{Z} = \{\mathbf{z}_i\}_{i=1}^N$  is zero, that is,  $\boldsymbol{\mu}_z = \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i = \mathbf{0}$ .
2. The covariance of  $\mathcal{Z}$  is diagonal, that is  $\boldsymbol{\Sigma}_z$  is diagonal.
3.  $\|\mathbf{x}_i - \mathbf{x}_j\|_2 = \|\mathbf{z}_i - \mathbf{z}_j\|_2$  for all i and j .

---

1.

$$\boldsymbol{\mu}_z = \frac{1}{N} \sum_{i=1}^N \mathbf{z}_i = \frac{1}{N} \sum_{i=1}^N \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) = \frac{1}{N} \mathbf{U}^T \left( \sum_{i=1}^N \mathbf{x}_i - \sum_{i=1}^N \boldsymbol{\mu}_x \right) = \mathbf{U}^T (\boldsymbol{\mu}_x - \boldsymbol{\mu}_x) = \mathbf{0}$$

2. Lets have a look at  $\Sigma_z$ :

$$\begin{aligned}
\Sigma_z &= \mathbb{E} [ZZ^T] = \mathbb{E} [U^T (X - \mu_X) (X - \mu_X)^T U] \\
&= U^T \mathbb{E} [(X - \mu_X)(X - \mu_X)^T] U \\
&= U^T \Sigma_X U \\
&= \begin{pmatrix} \lambda_1 & \cdots & 0 \\ 0 & & \vdots \\ \vdots & \cdots & \lambda_d \end{pmatrix} U^T U \\
&= \begin{pmatrix} \lambda_1 & \cdots & 0 \\ 0 & & \vdots \\ \vdots & \cdots & \lambda_d \end{pmatrix} I = \begin{pmatrix} \lambda_1 & \cdots & 0 \\ 0 & & \vdots \\ \vdots & \cdots & \lambda_d \end{pmatrix}
\end{aligned}$$

3. Lets take a look at  $\|z_i - z_j\|_2$ :

$$\begin{aligned}
\|z_i - z_j\| &= \|U^T(x_i - \mu_x) - U^T(x_j - \mu_x)\|_2 \\
&= \|U^T x_i - U^T \mu_x - U^T x_j + U^T \mu_x\|_2 \\
&= \|U^T x_i - U^T x_j\|_2 = \|U^T (x_i - x_j)\|_2 \\
U^T &\text{ is orthonormal (this is a rigid transformation) and therefore} \\
&= \|x_i - x_j\|_2
\end{aligned}$$

■

## Geometric PCA

Let  $U_d \in \mathbb{R}^{D \times d}$  be a full rank matrix (with  $d \leq D$ ).

### 1.2.2

Show that exists an invertible matrix  $M \in \mathbb{R}^{d \times d}$  such that  $O = U_d M \in \mathbb{R}^{D \times d}$  is semi-orthogonal, that is:

$$O^T O = I_d$$

Let  $U_d \in \mathbb{R}^{D \times d}$  be a fully ranked matrix (with  $d \leq D$ ). Since  $U$  is fully ranked, we can use the compact SVD decomposition and denote  $U_d = A \Sigma_d B^T$  where  $A \in \mathbb{R}^{D \times d}$ ,  $B \in \mathbb{R}^{d \times d}$  are two orthogonal matrices and  $\Sigma_d \in \mathbb{R}^{d \times d}$  is a diagonal matrix. Now, using compact-SVD, we will examine only the top  $d$  rows/cols:

$$U_d = A_d \Sigma_d B^T \implies U_d B = A_d \Sigma_d \implies A_d = U_d B \Sigma_d^{-1}$$

We will notice that  $A_d$  is semi-orthogonal,  $A_d^T A_d = I_d$ . Denoting  $M = B \Sigma_d^{-1}$  we will have exactly what we wished for.

■

### 1.2.3

In order to prove that both problems have the same solution, we will examine  $\frac{1}{N}\|X - U_D U_D^T X\|_F^2$ .

$$\begin{aligned}
\frac{1}{N}\|X - U_d U_d^T X\|_F^2 &= \frac{\text{Tr}[(X - U_d U_d^T X)(X - U_d U_d^T X)^T]}{N} \\
&= \frac{\text{Tr}[(X - U_d U_d^T X)(X^T - X^T U_d U_d^T)]}{N} \\
&= \frac{\text{Tr}(X X^T - X X^T U_d U_d^T - U_d U_d^T X X^T + U_d U_d^T X X^T U_d U_d^T)}{N} \\
&= \frac{\text{Tr}(X X^T) - \text{Tr}(X X^T U_d U_d^T) - \text{Tr}(U_d U_d^T X X^T) + \text{Tr}(U_d U_d^T X X^T U_d U_d^T)}{N} \\
&= \frac{\text{Tr}(X X^T) - \text{Tr}(U_d^T X X^T U_d) - \text{Tr}(U_d^T X X^T U_d) + \text{Tr}(U_d^T U_d U_d^T X X^T U_d)}{N} \\
&= \frac{\text{Tr}(X X^T) - \text{Tr}(U_d^T X X^T U_d) - \text{Tr}(U_d^T X X^T U_d) + \text{Tr}(U_d^T X X^T U_d)}{N} \\
&= \frac{1}{N} \text{Tr}(X X^T) - \frac{1}{N} \text{Tr}(U_d^T X X^T U_d) \\
&= \text{Tr}(\Sigma_X) - \frac{1}{N} \text{Tr}(U_d^T X X^T U_d)
\end{aligned}$$

Since  $\text{Tr}(\Sigma_X)$  and  $\frac{1}{N}$  is constant, then  $\|X - U_D U_D^T X\|_F^2$  is maximal when  $\text{Tr}(U_d^T X X^T U_d)$  is minimal. ■

### 1.2.4

We will denote  $\epsilon = X - \hat{X}$ , where  $\epsilon = \{\epsilon_i\}$  is the matrix of the errors,  $\epsilon_i = x_i - \hat{x}_i$ .

We need to prove that  $\text{Tr}(\Sigma_\epsilon) = \frac{1}{N}\|X - U_d U_d^T X\|_F^2$

$$\begin{aligned}
\text{Tr}(\Sigma_\epsilon) &= \text{Tr}(\mathbb{E}[\epsilon \epsilon^T]) \\
&= \text{Tr}(\mathbb{E}[(X - \hat{X})(X - \hat{X})^T]) \\
&= \text{Tr}(\mathbb{E}[(X - U_d U_d^T X)(X - U_d U_d^T X)^T]) \\
&= \text{Tr}(\frac{1}{N}[(X - U_d U_d^T X)(X - U_d U_d^T X)^T]) \\
&= \frac{1}{N} \text{Tr}([(X - U_d U_d^T X)(X - U_d U_d^T X)^T]) \\
&= \frac{1}{N} \|X - U_d U_d^T X\|_F^2
\end{aligned}$$

We notice that  $\text{Tr}(U_d^T X X^T U_d) = \text{Tr}(\Sigma_Z)$

Because of 1.1.3:

$$\frac{1}{N}\|X - U_d U_d^T X\|_F^2 = \text{Tr}(\Sigma_X) - \frac{1}{N} \text{Tr}(U_d^T X X^T U_d)$$

$$\text{Tr}(\Sigma_\epsilon) = \text{Tr}(\Sigma_X) - \text{Tr}(\Sigma_Z) \quad \blacksquare$$

### 1.2.5

Lets assume  $U_d \in \mathbb{R}^{D \times d}$  be the top d eigenvalues corresponding to d largest eigenvalues of  $\Sigma_x$

We need to prove that  $\text{Tr}(\Sigma_\epsilon) = \sum_{i=d+1}^D \lambda_i(\Sigma_x)$

we know that  $\text{Tr}(\Sigma_x) = \sum_{i=1}^D \lambda_i(X)$

Lets define  $Z = U^T X$

We know that  $\text{Tr}(\Sigma_Z) = \text{Tr}(U_d^T \Sigma_x U_d) = \text{Tr}(U_d U_d^T \Sigma_x) = \text{Tr}(I_d \Sigma_x) = \sum_{i=1}^d \lambda_i(X)$

from 1.2.4, we know that :  $\text{Tr}(\Sigma_\epsilon) = \text{Tr}(\Sigma_X) - \text{Tr}(\Sigma_Z)$

if we combine all of the data together we get that:

$$\text{Tr}(\Sigma_\epsilon) = \text{Tr}(\Sigma_X) - \text{Tr}(\Sigma_Z) = \sum_{i=1}^D \lambda_i(X) - \sum_{i=1}^d \lambda_i(X) = \sum_{i=d+1}^D \lambda_i(X) \quad \blacksquare$$


---

## High-dimensional data PCA

Consider the data  $\mathbf{X} \in \mathbb{R}^{D \times N}$  where  $D > N$ .

### 1.2.6

- Provide a (tight) upper bound on the number of non-zero eigenvalues.
  - Consequently, can you apply PCA to  $\mathbf{X} \in \mathbb{R}^{D \times N}$  to obtain  $\mathbf{Z} \in \mathbb{R}^{d \times N}$  with  $d < D$  such that there is no loss of information?  
Explain your answer.
- 

- We will notice that the number of non zero eigenvalues would is always bound by the matrix rank, i.e  $\text{rank}(X) \leq N$  and therefore,  $N$  is the (tight) upper bound.
  - Yes, its possible. As we seen in class, the top d eigenvectors of  $\Sigma_x$  are the same as the top (left) singular vectors of  $X$ . Therefore, if  $d \geq N$  we will not loss any information. \blacksquare
- 

## Rank minimization

- Let  $\mathbf{A} \in \mathbb{R}^{D \times N}$ .
- Consider the following rank minimization problem:

$$\begin{cases} \min_{\mathbf{M} \in \mathbb{R}^{D \times N}} \|\mathbf{A} - \mathbf{M}\|_F^2 \\ \text{s.t. } \text{rank}(\mathbf{M}) \leq d \end{cases}$$

### 1.2.7

- Solve the optimization problem.
  - Write your final solution using the (truncated) matrices obtained by the SVD decomposition of  $\mathbf{A}$ , namely,  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$
- 

We will notice, that because  $\text{rank}(M) = d$  then  $d \leq \min(D, N)$  and therefore we have two matrices  $B \in \mathbb{R}^{D \times d}, C \in \mathbb{R}^{d \times N}$  with ranks  $d$  such that  $M = BC$ . So we actually have

$$\begin{aligned} \min_{\mathbf{M} \in \mathbb{R}^{D \times N}} \|\mathbf{A} - \mathbf{M}\|_F^2 &= \min_{\mathbf{B} \in \mathbb{R}^{D \times d}, \mathbf{C} \in \mathbb{R}^{d \times N}} \|\mathbf{A} - \mathbf{BC}\|_F^2 = \\ &= \min_{\mathbf{U}_d \in \mathbb{R}^{D \times d}, \mathbf{Z} \in \mathbb{R}^{d \times N}} \|\mathbf{A} - \mathbf{U}_d \mathbf{Z}\|_F^2 \end{aligned}$$

This is equivalent to a PCA problem which we saw. Therefore the optimal solution is for

$$\begin{aligned} Z &= U_d^T A \\ U_d^T U_d &= I_d \end{aligned}$$

and when  $A = U\Sigma V^T$  the SVD decomposition,  $M = U_d \Sigma_d V_d^T$ .

## 2 KPCA

### 2.1

#### 2.1.1

There is  $J = I - \frac{1}{N}11^T \in \mathbb{R}^{N \times N}$   
 We need to show that  $J^2 = J$

$$\begin{aligned} J^2 &= J * J = (I - \frac{1}{N}11^T)(I - \frac{1}{N}11^T) \\ &= I - \frac{1}{N}11^T - \frac{1}{N}11^T + \frac{1}{N^2}11^T11^T \\ &=^* I - \frac{1}{N}11^T - \frac{1}{N}11^T + \frac{1}{N^2}N11^T \\ &= I - \frac{1}{N}11^T - \frac{1}{N}11^T + \frac{1}{N}11^T = I - \frac{1}{N}11^T = J \end{aligned}$$

where in  $*$ , we know that:

$$11^T11^T = \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix} \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix} = \begin{pmatrix} N & \cdots & N \\ \vdots & & \vdots \\ N & \cdots & N \end{pmatrix} = N \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix} = N11^T \quad \blacksquare$$

#### 2.1.2

Let  $X \in \mathbb{R}^{D \times N}$ ,  $\Sigma_x = XX^T$ ,  $K_x = X^T X$  and we know that  $\Sigma_x u_i = \lambda_i u_i$  so  $\lambda_i u_i = \Sigma_x u_i = XX^T u_i$ .  
 We will multiply both sides with  $X^T$  and we will get

$$\begin{aligned} X^T \lambda_i u_i &= \lambda_i X^T u_i = X^T X X^T u_i \\ &\implies (X^T X) X^T u_i = \lambda_i X^T u_i \\ &\implies K_x \cdot (X^T u_i) = \lambda_i \cdot (X^T u_i) \end{aligned}$$

Therefore  $\lambda_i$  is an eigenvalue of  $K_x$  and the corresponding eigenvector is  $X^T u_i$ . ■

---

## Kernel functions

Let  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  and consider  $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^N$ .

#### 2.1.3

Show that if  $k$  can be written as an inner product, that is

$$k(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

for some  $\phi$ , then, the matrix defined by:

$$\mathbf{K}_x[i, j] = k(\mathbf{x}_i, \mathbf{x}_j)$$

is an SPSP matrix, namely,  $\mathbf{K}_x \succeq 0$ .

---

### 2.1.3

Lets start by showing that every matrix multiplication in the shape of  $A^T A$  is SPSD.

In order to do so, we need to show that  $v^T A^T A v \geq 0$

$$v^T A^T A v = (Av)^T Av = \langle Av, Av \rangle = \|Av\|^2 \geq 0.$$

And now only thing left to show show is that  $K_x$  can be written as  $K_x = \phi^T \phi$ .

Let  $\Phi$  be a matrix of applications of  $\phi$  on every instance  $x$ , then  $K[i, j]$ :

$$K_x[i, j] = \langle \phi(x_i), \phi(x_j) \rangle = \Phi^T \Phi[i, j]$$

Which means that every element in  $K_x$  equals to every element in  $\Phi^T \Phi$ , and therefore  $K_x = \Phi^T \Phi$ .  
Therefore  $K$  is SPSD. ■

---

#### 2.1.4.1

Let  $A$  be an SPD matrix, and let:

$$k(x_i, x_j) = x_i^T A x_j$$

Prove or disprove:  $k$  is a kernel function.

We will prove that  $k(x_i, x_j) = x_i^T A x_j$  is a kernel function.

Since  $A$  is SPD, we can use EVD:  $A = U \Lambda U^T$ , where  $U \in \mathbb{R}^{d \times d}$  Lets write

$$\phi(x) = \begin{bmatrix} \sqrt{\lambda_1} \langle x, u_1 \rangle \\ \sqrt{\lambda_2} \langle x, u_2 \rangle \\ \vdots \\ \sqrt{\lambda_d} \langle x, u_d \rangle \end{bmatrix}$$

Now, lets show that at  $k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$

$$\begin{aligned} k(x_i, x_j) &= x_i^T A x_j \\ &= x_i^T U \Lambda U^T x_j \\ &= \begin{bmatrix} \langle x_i, u_1 \rangle \\ \langle x_i, u_2 \rangle \\ \vdots \\ \langle x_i, u_d \rangle \end{bmatrix} \Lambda \begin{bmatrix} \langle u_1, x_j \rangle \\ \langle u_2, x_j \rangle \\ \vdots \\ \langle u_d, x_j \rangle \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 \langle x_i, u_1 \rangle \\ \lambda_2 \langle x_i, u_2 \rangle \\ \vdots \\ \lambda_d \langle x_i, u_d \rangle \end{bmatrix} \begin{bmatrix} \langle u_1, x_j \rangle \\ \langle u_2, x_j \rangle \\ \vdots \\ \langle u_d, x_j \rangle \end{bmatrix} \\ &= \lambda_1 \langle x_i, u_1 \rangle \langle u_1, x_j \rangle + \dots + \lambda_d \langle x_i, u_d \rangle \langle u_d, x_j \rangle \\ &= \sqrt{\lambda_1} \langle x_i, u_1 \rangle \sqrt{\lambda_1} \langle u_1, x_j \rangle + \dots + \sqrt{\lambda_d} \langle x_i, u_d \rangle \sqrt{\lambda_d} \langle u_d, x_j \rangle \\ &= \begin{bmatrix} \sqrt{\lambda_1} \langle x_i, u_1 \rangle \\ \sqrt{\lambda_2} \langle x_i, u_2 \rangle \\ \vdots \\ \sqrt{\lambda_d} \langle x_i, u_d \rangle \end{bmatrix} \begin{bmatrix} \sqrt{\lambda_1} \langle u_1, x_j \rangle \\ \sqrt{\lambda_2} \langle u_2, x_j \rangle \\ \vdots \\ \sqrt{\lambda_d} \langle u_d, x_j \rangle \end{bmatrix} \\ &= \langle \phi(x_i), \phi(x_j) \rangle \end{aligned}$$



### 2.1.4.2

Let  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$  and consider:

$$k(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$$

Prove or disprove:  $k$  is a kernel function.

We will prove that  $k(x_i, x_j) = (1 + x_i^T x_j)^2$  is a kernel function.

Lets denote  $\phi(x) : \mathbb{R}^d \rightarrow \mathbb{R}^{2d+1}$  :

$$\phi(x) = [1 \quad \sqrt{2}x_1 \quad \cdots \sqrt{2}x_d \quad x_1^2 \quad \cdots x_d^2]^T$$

Now, lets show that  $k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$

$$\begin{aligned} k(\mathbf{x}_i, \mathbf{x}_j) &= (1 + \mathbf{x}_i^T \mathbf{x}_j)^2 \\ &= (1 + \mathbf{x}_i^T \mathbf{x}_j)(1 + \mathbf{x}_i^T \mathbf{x}_j) \\ &= 1 + 2\mathbf{x}_i^T \mathbf{x}_j + (\mathbf{x}_i^T \mathbf{x}_j)^2 \\ &= 1 + \left\langle \sqrt{2}\mathbf{x}_i, \sqrt{2}\mathbf{x}_j \right\rangle + (\langle \mathbf{x}_i, \mathbf{x}_j \rangle)^2 \\ &= \left\langle [1 \quad \sqrt{2}x_{i1} \quad \cdots \sqrt{2}x_{id} \quad x_{i1}^2 \quad \cdots x_{id}^2]^T, [1 \quad \sqrt{2}x_{j1} \quad \cdots \sqrt{2}x_{jd} \quad x_{j1}^2 \quad \cdots x_{jd}^2]^T \right\rangle \\ &= \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \end{aligned}$$

- Consider  $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^N$ , and consider the kernel:

$$k(\mathbf{x}_i, \mathbf{x}_j) := \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

for some  $\phi$ .

- Let:

$$\tilde{k}(\mathbf{x}_i, \mathbf{x}_j) := \langle \phi(\mathbf{x}_i) - \boldsymbol{\mu}_\phi, \phi(\mathbf{x}_j) - \boldsymbol{\mu}_\phi \rangle$$

be the centered version, where:

$$\boldsymbol{\mu}_\phi = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}_i)$$

### 2.1.5

Show that  $\tilde{k}$  can be written using only  $k$ , and without using  $\phi$  and  $\boldsymbol{\mu}_\phi$  explicitly.

$$\begin{aligned} \hat{k}(x_i, x_j) &= \langle \phi(x_i) - \boldsymbol{\mu}_\phi, \phi(x_j) - \boldsymbol{\mu}_\phi \rangle \\ &= (\phi(x_i) - \boldsymbol{\mu}_\phi)(\phi(x_j) - \boldsymbol{\mu}_\phi)^T \\ &= (\phi(x_i) - \boldsymbol{\mu}_\phi)(\phi(x_j)^T - \boldsymbol{\mu}_\phi^T) \\ &= \phi(x_i)\phi(x_j)^T - \phi(x_i)\boldsymbol{\mu}_\phi^T - \boldsymbol{\mu}_\phi\phi(x_j)^T + \boldsymbol{\mu}_\phi\boldsymbol{\mu}_\phi^T \\ &= k(x_i, x_j) - \frac{1}{N} \sum_{t=1}^N \phi(x_i)\phi(x_t)^T - \sum_{t=1}^N \phi(x_t)\phi(x_j)^T + \frac{1}{N^2} \sum_{t_1=1}^N \sum_{t_2=1}^N \phi(x_{t_1})\phi(x_{t_2})^T \\ &= k(x_i, x_j) - \frac{1}{N} \sum_{t=1}^N k(x_i, x_t) - \sum_{t=1}^N k(x_t, x_j) + \frac{1}{N^2} \sum_{t_1=1}^N \sum_{t_2=1}^N k(x_{t_1}, x_{t_2}) \end{aligned}$$

■

### 2.1.6

Prove or disprove:  $\widetilde{K}_x$  is an SPD matrix.

Let  $K_x \in \mathbb{R}^{N \times N}$  be a kernel matrix for some kernel function  $k$ ,  $K[i, j] = k(x_i, x_j)$

Let  $\hat{K}_x$  be the centered version, that is  $\hat{K}_x = JK_xJ$ , where  $J = I - \frac{1}{N}11^T$ .

We will prove that  $\hat{K}_x$  is not an SPD, by showing a counter example:

Let  $K_x = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$  and  $J = \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{pmatrix}$ . Centring  $K_x$  will give us  $\hat{K}_x = \begin{pmatrix} \frac{1}{4} & -\frac{1}{4} \\ -\frac{1}{4} & \frac{1}{4} \end{pmatrix}$ .

Now, let's consider the vector  $v = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  for which  $v^T \hat{K}_x v = 0$ , and therefore  $\hat{K}_x$  is not an SPD. ■

## Out Of sample extension

- Let  $K_x$  be the kernel matrix obtained from the training set  $\mathcal{X} = \{\mathbf{x}_i \in \mathbb{R}^D\}_{i=1}^N$ .
- Let  $Z \in \mathbb{R}^{d \times N}$  be the low-dimensional representation obtained by applying KPCA, that is:

$$Z = \Sigma_d V_d^T$$

where  $V \Sigma V^T = J K_x J$  is an eigendecomposition (see lecture notes).

- Let  $X^* \in \mathbb{R}^{D \times N^*}$  be a set of new unseen data-points.

### 2.1.7

Write an expression (in a matrix form) for  $Z^* \in \mathbb{R}^{d \times N^*}$ , the KPCA out of sample extension applied to  $X^*$ .

We will denote  $\Phi^* = \phi(X^*)$  the application of  $\Phi$  on all the unseen data-points  $X^*$ . We will also denote  $K_x^* = \Phi^T \Phi^*$  and  $\widetilde{K}_x^* = \widetilde{\Phi}^T \widetilde{\Phi}^*$ .

Then from what we learned in the lecture, we can represent the KPCA of the OOS extension applied on  $X^*$  as:

$$Z^* = \Sigma_d^{-1} V_d^T \widetilde{\Phi}^T \widetilde{\Phi}^* = \Sigma_d^{-1} V_d^T \widetilde{K}_x^*$$

Since  $\widetilde{K}_x^* = J(K_x^* - \frac{1}{N} K_x 1_N 1_D^T)$

$$Z^* = \Sigma_d^{-1} V_d^T \widetilde{K}_x^* = \Sigma_d^{-1} V_d^T J(K_x^* - \frac{1}{N} K_x 1_N 1_D^T)$$

■

### 2.1.8

Let  $\chi^* = \{x_i^*\}_{i=1}^N$  be a subset of the training set  $\chi$ , let  $X^* \in \mathbb{R}^{D \times N^*}$  be the matrix form of  $\chi^*$  and let  $Z^* \in \mathbb{R}^{d \times N}$  be the low dimension representation obtained by the training encoding.

We need to prove that the out of sampling encoding applied to  $X^*$  coincide with the training encoding  $Z^*$ .

We know that  $X V_d \Sigma_d^{-1} = U_d$ , that  $K = \phi(X) \phi(X)^T = V \Sigma^2 V^T$ , and that  $z^* = U_d^T \phi(x^*) = \Sigma_d^{-1} V_d^T \phi(X)^T \phi^*(x)$

Lets assume that the out of sample  $X^* = X$

$$\begin{aligned}
Z^* &= \Sigma_d^{-1} V_d^T \tilde{\Phi}(X)^T \tilde{\Phi}(x^*) \\
&= \Sigma_d^{-1} V_d^T \tilde{\Phi}(X)^T \tilde{\Phi}(X) \\
&= \Sigma_d^{-1} V_d^T \tilde{K} \\
&= \Sigma_d^{-1} V_d^T V \Sigma^2 V^T \\
&= \Sigma_d V_d^T = Z
\end{aligned}$$

Lets assume that the out of sample contains only one sample  $x_i$ , and lets take a look at  $z_i$ :

$$\begin{aligned}
Z^* &= \Sigma_d^{-1} V_d^T \tilde{\Phi}(X)^T \tilde{\phi}(x_i) \\
&= \Sigma_d^{-1} V_d^T \tilde{\Phi}(K_x)_i \\
&= [\Sigma_d^{-1} V_d^T \tilde{\Phi}(K_x)]_i \\
&= (\Sigma_d^{-1} V_d^T V \Sigma^2 V^T)_i \\
&= (\Sigma_d V_d^T) = Z_i
\end{aligned}$$

So for each out of sample  $x_i$ , the encoding is  $Z_i$ . ■