

ADVANCED MACHINE LEARNING

Kernel PCA



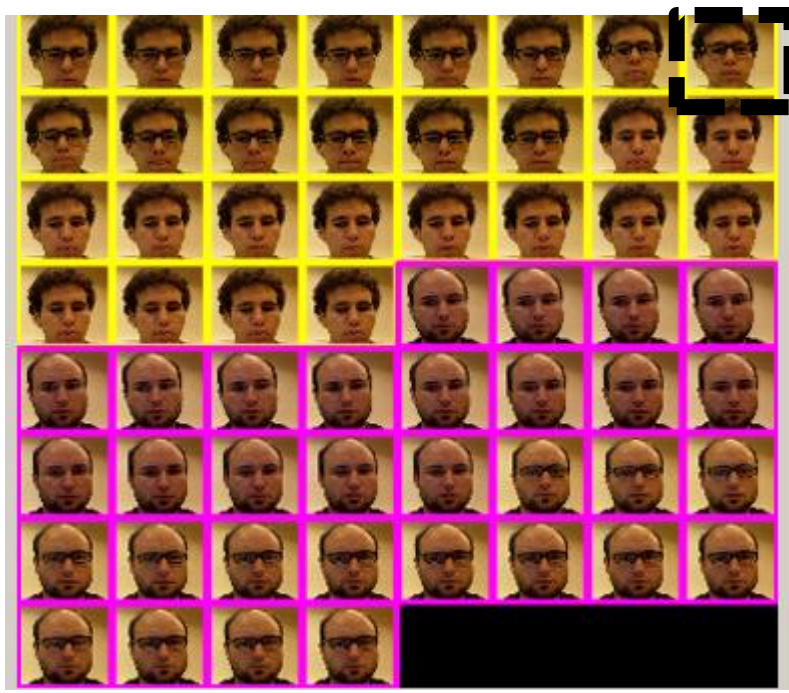
ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Overview Today's Lecture

- Brief Recap of Classical Principal Component Analysis (PCA)
- Derivation of kernel PCA
- Exercises to develop a geometrical intuition of kernel PCA

Principal Component Analysis: Overview

Take samples of two classes (yellow and pink classes)



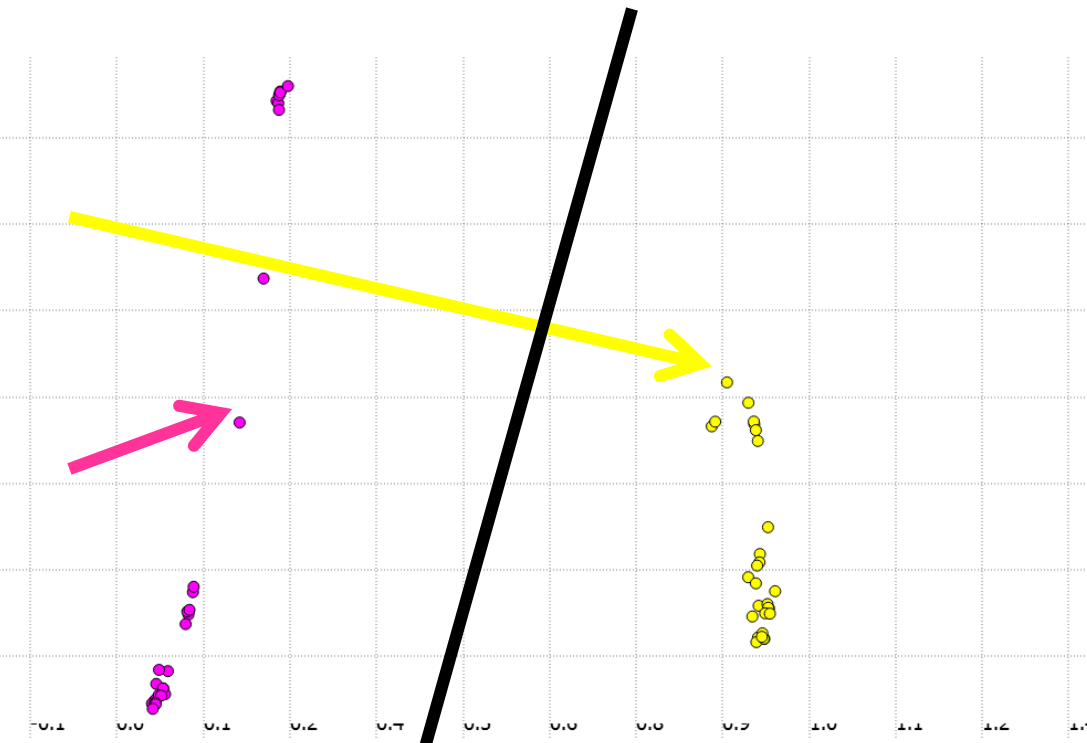
Each image is a high-dimensional vector

$$x \in \mathbb{R}^{320 \cdot 240 \cdot 3 = 230400}$$

Principal Component Analysis: Overview

Project the images onto a lower dimensional space

$$y \in \mathbb{R}^2 \text{ through matrix } A \in \mathbb{R}^{2 \times 230400} : y = Ax$$



Separating Line

Principal Component Analysis: Overview

Project the images onto a lower dimensional space

$y \in \mathbb{R}^2$ through matrix $A \in \mathbb{R}^{2 \times 230400}$: $y = Ax$

**What is A?
PCA discovers the matrix A**

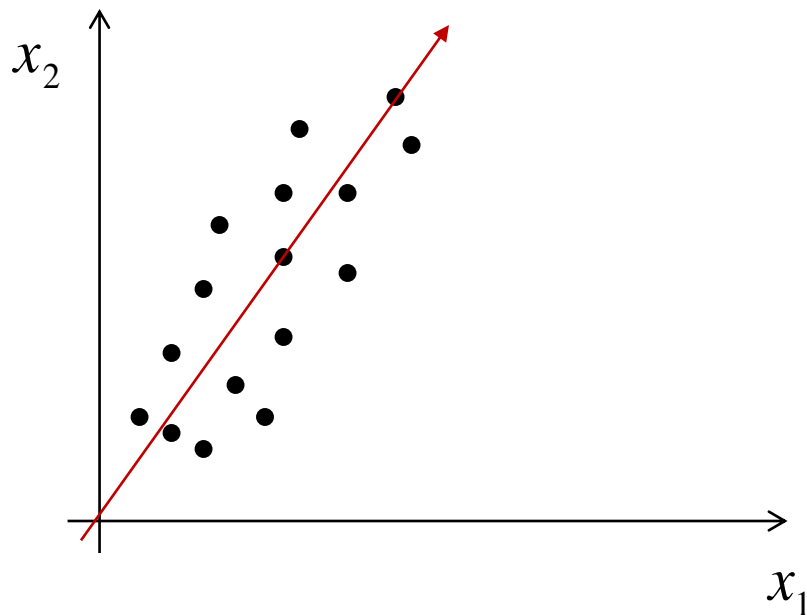


Principal Component Analysis: Overview

Infinite number of choices for projection matrix A

→ need criteria to reduce the choice

1: minimum information loss (minimal reconstruction error)



What is the 2D to 1D projection that minimizes the reconstruction error?

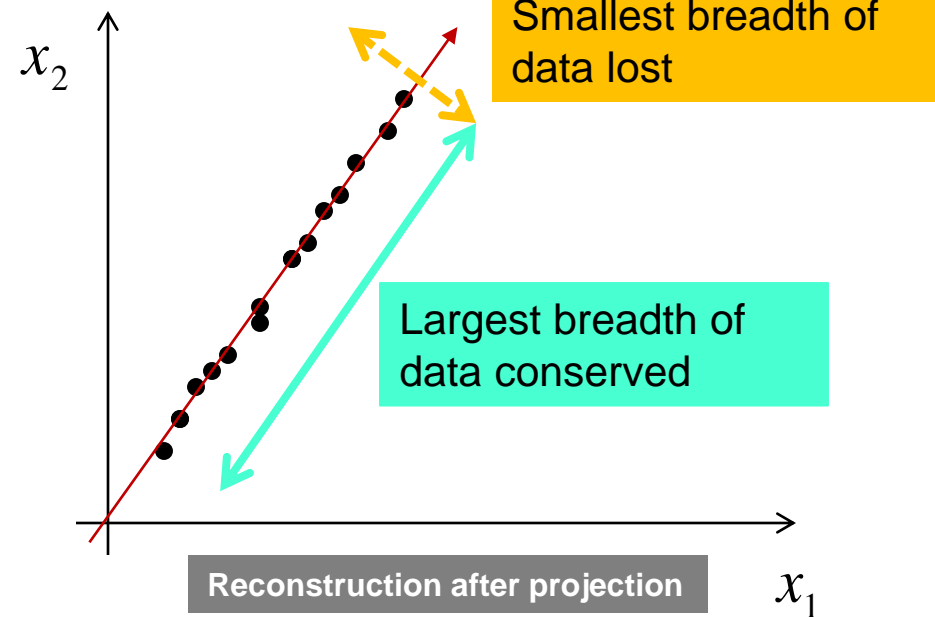
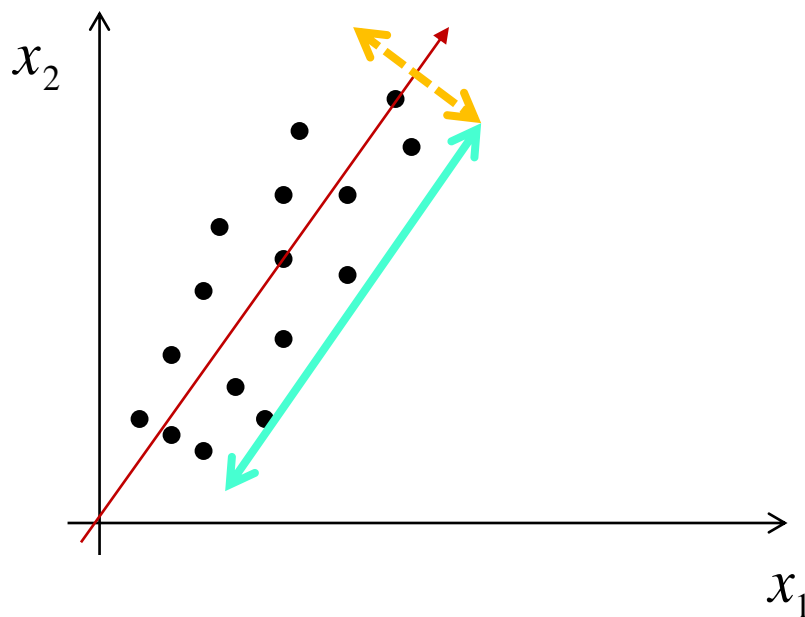
Principal Component Analysis: Overview

Infinite number of choices for projection matrix A

→ need criteria to reduce the choice

1: minimum information loss (minimal reconstruction error)

2: equivalent to finding the direction with maximum variance



What is the 2D to 1D projection that minimizes the reconstruction error?

Principal Component Analysis: Overview

Dataset $X = [x^1 \ x^2 \ \dots \ x^M]$ (data is centered, i.e. $E\{X\} = 0$)

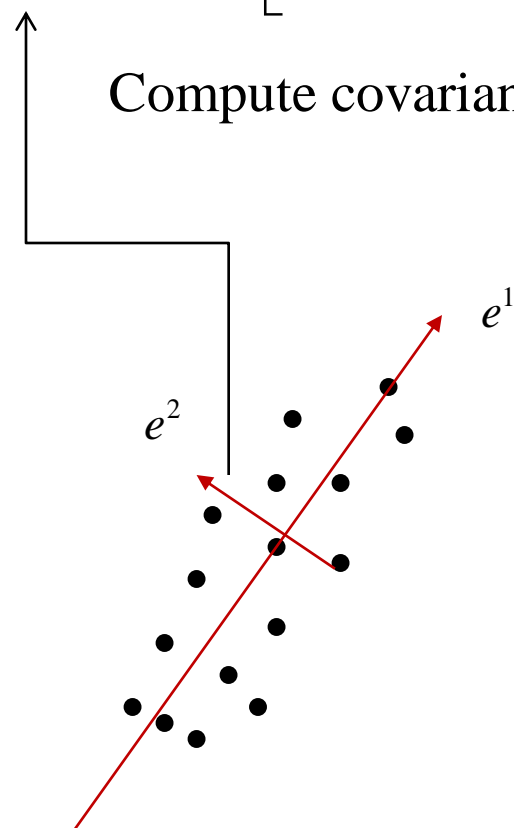
Compute covariance matrix of dataset X : $C = \frac{1}{M} E\{XX^T\}$

Find eigenvalue decomposition: $C = V\Lambda V^T$

$V = [e^1 \ \dots \ e^N]$: matrix of eigenvectors

Λ : Diagonal matrix of eigenvalues

Order $e^1 \ \dots \ e^N$, s.t. $\lambda_1 \geq \lambda_2 \dots \geq \lambda_N$



The eigenvectors form a basis of the space.

e^1 is aligned with the axis of maximum variance.

Project data onto eigenvectors.

Remove projections with low λ (noise).

PCA for Data Compression

Original image is encoded in $x \in \mathbb{R}^N$. Compressed image is $y \in \mathbb{R}^{p=0.1N}$

$$y = A_p x,$$

Rows of A_p contains 1st p eigenvectors

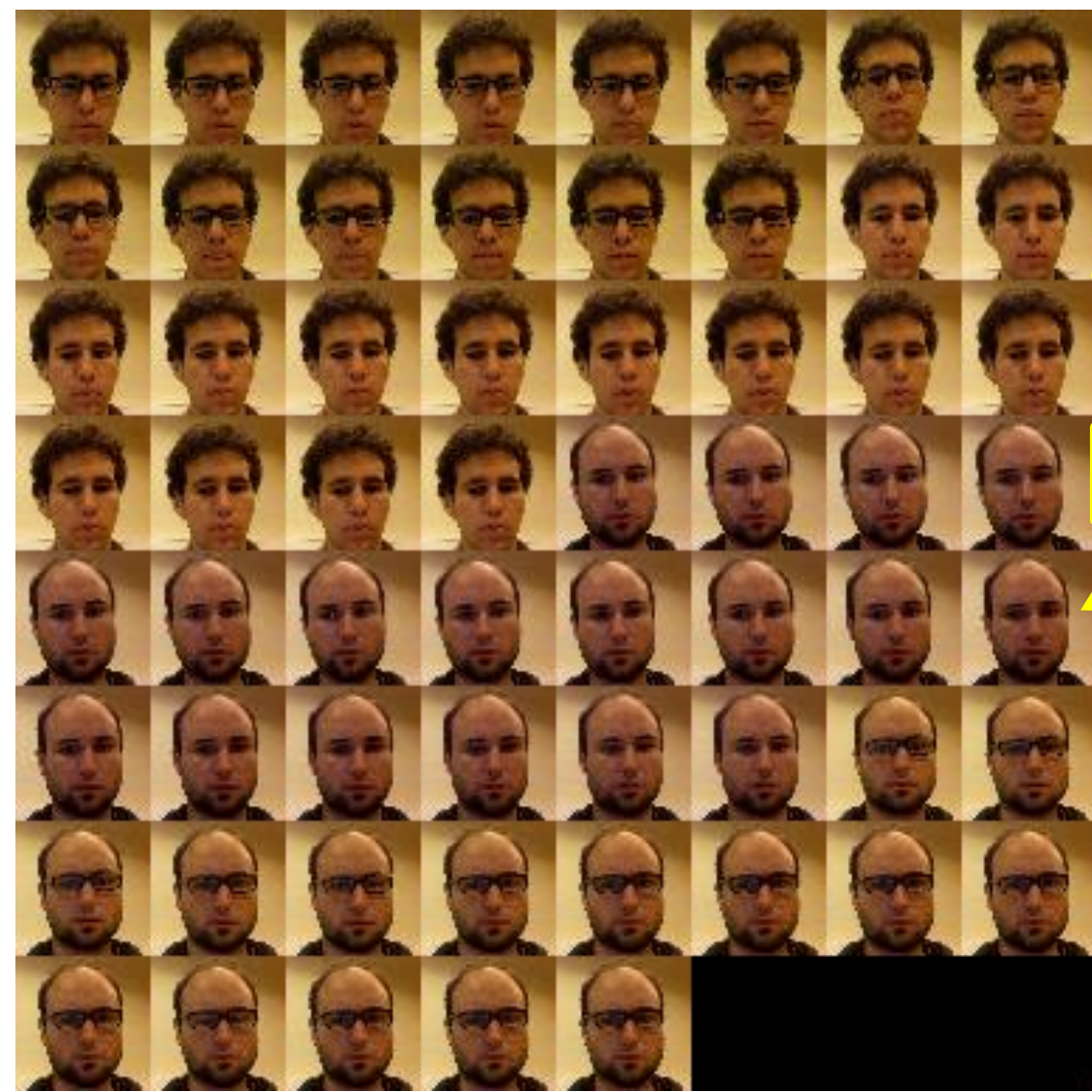


Original Image



Image compressed 90%

PCA for Feature Extraction



Results of decomposition with Principal Component Analysis: eigenvectors



Encapsulate main differences across groups of images (in the first eigenvectors)

Detailed features (glasses) get encapsulated next (in the following eigenvectors)

Principal Component Analysis: Pros & Cons

Advantages:

- a) The projection through PCA ensures minimal reconstruction error.
- b) The projection does not distort the space (rotation in space).
 - Ease of visualization/interpretation: The features that appear in the projections are often interpretable visually.

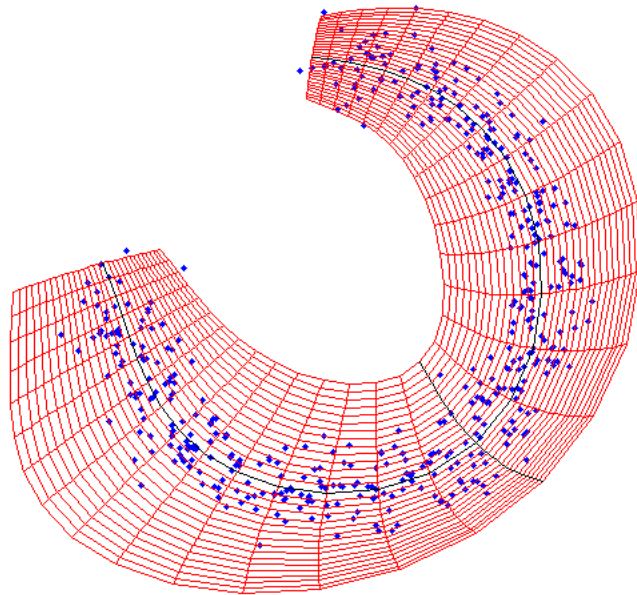
Limitations:

- a) PCA assumes a linear transformation:
 - With centering of data, one can only do a rotation in space.
- b) It fails at finding directions that require a non-linear transformation.

Revisiting the hypotheses of PCA

PCA assumed a *linear* transformation

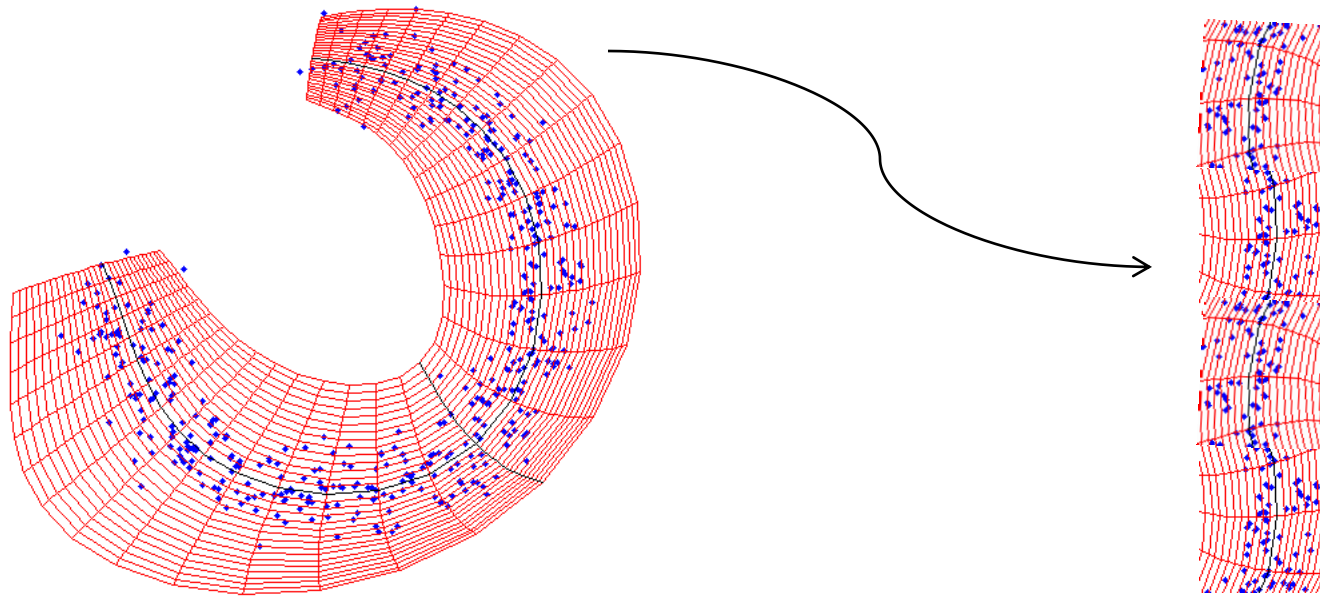
→ *Non-linear* PCA ([Kernel PCA](#)): find a non-linear embedding of the data and then perform linear PCA.



Recall: Principle of kernel Methods

Going back to linearity

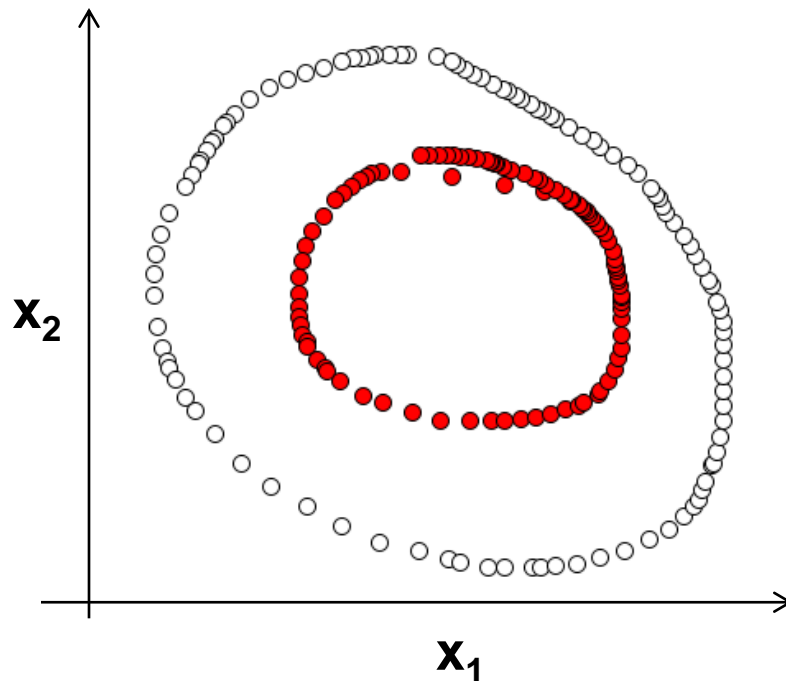
Find a non-linear transformation that send the data in a space where linear computation is again feasible.



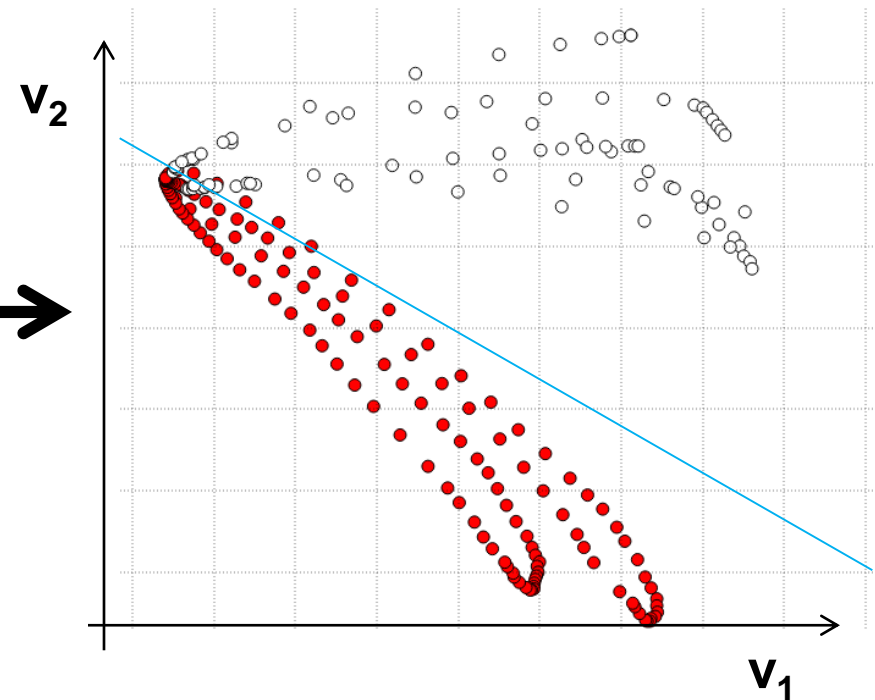
Kernel PCA: Principle

Determine a transformation which brings out features of the data so as to make subsequent computation easier.

Original Space



After Lifting Data in Feature Space



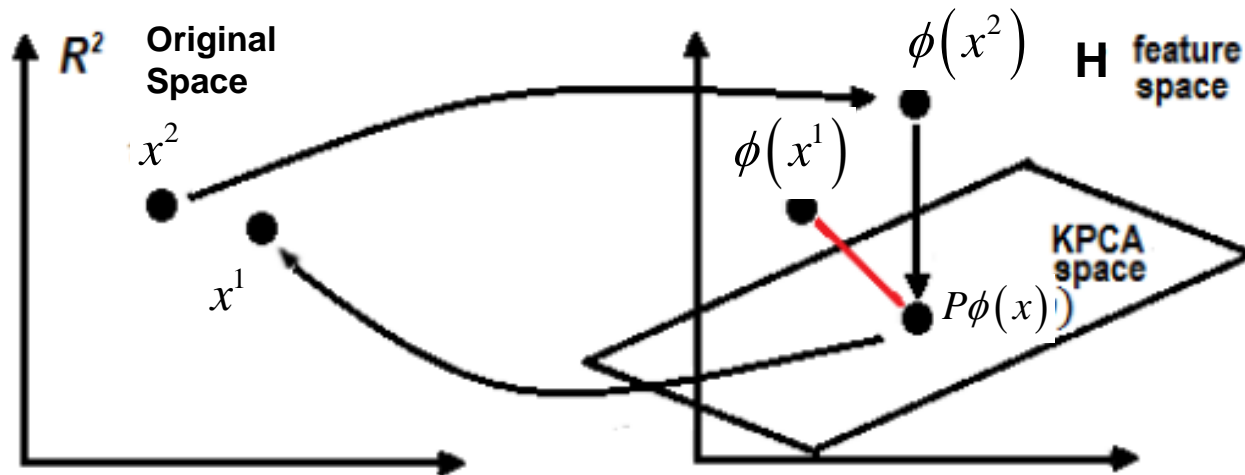
Example above: Data becomes *linearly* separable when using a rbf kernel and projecting onto first 2 PC-s of kernel PCA.

Kernel PCA: Principle

Idea: Send the data X into a **feature space H** through the **nonlinear map ϕ** .

$$X = \{x^i \in \mathbb{R}^N\}_{i=1\dots M} \mapsto \phi(X) = (\phi(x^1), \dots, \phi(x^M))$$

Perform **linear pca** in feature space and **project** into set of **eigenvectors** in feature space



Kernel PCA: Principle

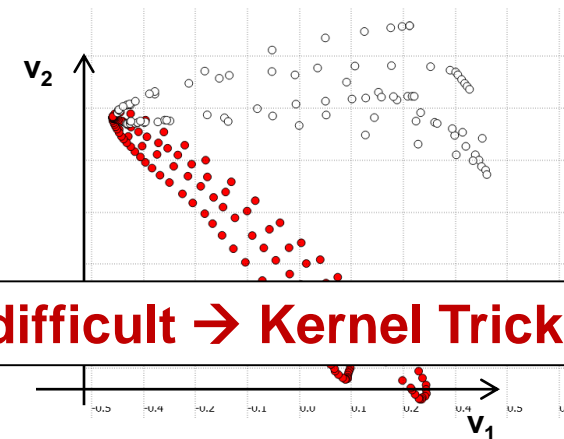
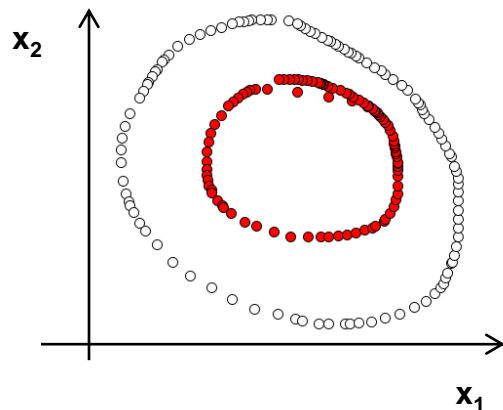
Idea: Send the data X into a **feature space H** through the **nonlinear map ϕ** .

$$X = \{x^i \in \mathbb{R}^N\}_{i=1 \dots M} \mapsto \phi(X) = (\phi(x^1), \dots, \phi(x^M))$$

Perform **linear pca** in feature space and **project** into set of **eigenvectors** in feature space

Data projected onto the two first principal components in feature space

Original Space



Determining ϕ is difficult \rightarrow Kernel Trick

Linear PCA in Feature Space

Sending the data in feature space through ϕ :

$$\phi : X \rightarrow H \quad x \mapsto \phi(x)$$

Assume that, in feature space H , the data are centered:

$$\sum_{i=1}^M \phi(x^i) = 0$$

The covariance matrix in the feature space is:

$$C_{\phi} = \frac{1}{M} FF^T$$

The columns $i = 1 \dots M$ of F are composed of $\phi(x^i)$.

Linear PCA in Feature Space

As in the original space, in feature space, the covariance matrix can be diagonalized and we have now to find the eigenvalues $\lambda_i \geq 0$, satisfying:

$$C_\phi v^i = \lambda_i v^i$$

Primal eigenvalue problem: Finding the eigenvectors v of C_ϕ

Not possible in feature space!

=> Formulate everything as a dot product and use kernel trick!

From Linear PCA to Kernel PCA

Each eigenvector v^1, \dots, v^M can be expressed as a linear combination of the images of the datapoints:

Rewriting PCA in terms of dot products:

Using $C_\phi v^i = \frac{1}{M} \sum_{j=1}^M \phi(x^j) \phi(x^j)^T v^i$ with $C_\phi v^i = \lambda_i v^i$

we obtain, $v^i = \frac{1}{\lambda_i M} \sum_{j=1}^M \phi(x^j) \underbrace{\phi(x^j)^T v^i}_{\alpha_j^i} = \frac{1}{\lambda_i M} \sum_{j=1}^M \alpha_j^i \phi(x^j).$

Scalar


Linear PCA in Feature Space

Multiplying the equation:

$$C_{\phi} v^i = \lambda_i v^i$$

by $\phi(x^j)^T$, on both sides, we have:

$$\langle \phi(x^j), C_{\phi} v^i \rangle = \lambda_i \langle \phi(x^j), v^i \rangle, \quad \forall i, j = 1, \dots, M$$



$$C_{\phi} v^i = \frac{1}{M} \sum_{k=1}^M \phi(x^k) \phi(x^k)^T v^i \quad \leftarrow \quad v^i = \frac{1}{\lambda_i M} \sum_{l=1}^M \alpha_l^i \phi(x^l)$$

Linear PCA in Feature Space

$$\Rightarrow \frac{1}{\lambda_i M^2} \sum_{k=1}^M \left\langle \underbrace{\phi(x^j), \phi(x^k)}_{\sum_k K_{jk}} \sum_{l=1}^M \alpha_l^i \underbrace{\langle \phi(x^k), \phi(x^l) \rangle}_{K_{kl}} \right\rangle = \frac{1}{M} \sum_{l=1}^M \alpha_l^i \underbrace{\langle \phi(x^j), \phi(x^l) \rangle}_{K_{jl}}$$

Use the kernel trick: $k(x^i, x^j) := K_{ij} = \langle \phi(x^i), \phi(x^j) \rangle$

→ Eigenvalue problem of the form:

$$K \alpha^i = M \lambda_i \alpha^i, \quad K : \text{Gram Matrix}$$

Dual eigenvalue problem: Finding the dual eigenvectors α^i .

Linear PCA in Feature Space

The solutions to the dual eigenvalue problem :
are given by all the eigenvectors $\alpha^1, \dots, \alpha^M$ with
non-zero eigenvalues $\lambda_1, \dots, \lambda_M$.

Kernel PCA finds at most M eigenvectors
 M : number of datapoints
 $M \gg N$ dimension of each datapoint

→ Eigenvalue problem of the form:

$$K \alpha^i = M \lambda_i \alpha^i, \quad K : \text{Gram Matrix}$$

Linear PCA in Feature Space

Request that the eigenvectors v of C_ϕ be normalized,

$$\text{i.e. } \langle v^i, v^i \rangle = 1 \quad \forall i = 1, \dots, M$$

is equivalent to asking that the dual eigenvectors $\alpha^1, \dots, \alpha^M$ are such that: $1 / \lambda^i = \|\alpha^i\|$.

Constructing the kPCA projections

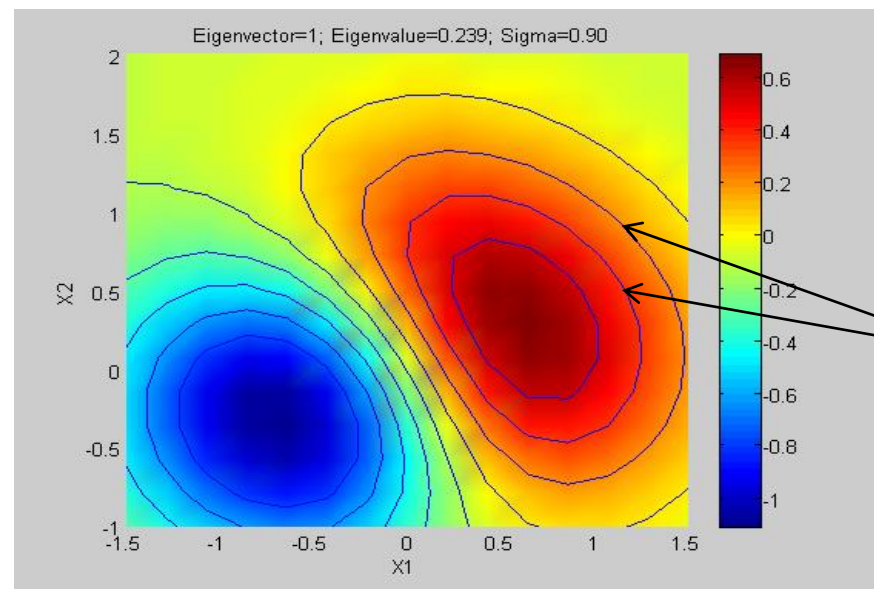
We cannot see the projection in feature space!

We can only compute the projections of each point onto each eigenvector

Projection of query point x onto eigenvector v^i :

$$\langle v^i, \phi(x) \rangle = \frac{1}{\lambda_i M} \sum_{j=1}^M \alpha_j^i \langle \phi(x^j), \phi(x) \rangle = \frac{1}{\lambda_i M} \sum_{j=1}^M \alpha_j^i k(x^j, x)$$

Sum over all training points



Isolines group points with equal projection:
 All points $x, s.t. \langle v^i, \phi(x) \rangle = cst.$

kPCA projections: Exercise

Recall: projection of query point x onto eigenvector v^i :

$$\langle v^i, \phi(x) \rangle = \frac{1}{\lambda_i M} \sum_{j=1}^M \alpha_j^i k(x^j, x)$$

where the α^i are the dual eigenvectors,
 solutions of the eigenvalue decomposition of K .

Consider a 2 – dimensional data – space, with two datapoints,
 and the RBF kernel:

$$k(x, x') = e^{-\frac{\|x-x'\|^2}{\sigma^2}}$$

- How many dual eigenvectors do you have and what is their dimension?
- Compute the eigenvectors and draw the isolines for the projections on each eigenvector.
- Repeat (b) for a homogeneous polynomial kernel with $p=2$: $k(x, x') = \langle x, x' \rangle^2$

kPCA projections: Exercise

Recall: projection of query point x onto eigenvector v^i :

$$\langle v^i, \phi(x) \rangle = \frac{1}{\lambda_i M} \sum_{j=1}^M \alpha_j^i k(x^j, x)$$

where the α^i are the dual eigenvectors,
 solutions of the eigenvalue decomposition of K .

Consider a 2 – dimensional data – space, with three equidistant datapoints,
 and the RBF kernel:

$$k(x, x') = e^{-\frac{\|x-x'\|^2}{\sigma^2}}$$

- How many dual eigenvectors do you have and what is their dimension?
- Compute the eigenvectors and draw the isolines for the projections on each eigenvector.
- Repeat (b) for a homogeneous polynomial kernel with $p=2$: $k(x, x') = \langle x, x' \rangle^2$
- What happens if you take 3 non-equidistant datapoints?

Curse of Dimensionality

Kernel PCA is very intensive computationally.

Computation of the eigenvectors requires eigenvalue decomposition of the Gram matrix (Kernel Matrix is $M \times M$) which grows quadratically with the number of data points M .

Computation of each projection in original space grows linearly with M too.

→ A variety of sparse methods have been proposed in the literature

Summary

- Kernel PCA offers an alternative to standard PCA to determine projections of the data while not assuming a linear transformation.
- It exploits the principle of the kernel trick, by replacing the inner product between pair of datapoints in the standard PCA (to compute the Covariance matrix) by the kernel function.
- The computation of kernel PCA is simple, yet it is expensive as it requires an eigenvalue decomposition of a very large matrix.
- Current research develops sparse versions of the algorithm.