PCA (Wrap-up) and Nonlinear Dimensionality Reduction via Kernel PCA

Piyush Rai

Machine Learning (CS771A)

Sept 7, 2016

Machine Learning (CS771A)

PCA (Wrap-up) and Nonlinear Dimensionality Reduction via Kernel PCA

・ロッ ・通ッ ・ヨッ ・ヨッ

Recap/Wrap-up of PCA

◆□▶ ◆□▶ ◆□▶ ◆□▶ ●□



イロン 不同と イヨン イヨン



• PCA basically does the following

(日) (同) (日) (日)



- PCA basically does the following
 - Learns the most important directions (new basis vectors) in the data

・ 同 ト ・ ヨ ト ・ ヨ ト



- PCA basically does the following
 - Learns the most important directions (new basis vectors) in the data
 - *Re*-represents data using the new basis vectors (change of basis)

(B)



- PCA basically does the following
 - Learns the most important directions (new basis vectors) in the data
 - Re-represents data using the new basis vectors (change of basis)
 - Throws away "uninteresting" directions without much loss of info. This gives a new lower-dimensional representation of the original data

メポト イヨト イヨト



- PCA basically does the following
 - Learns the most important directions (new basis vectors) in the data
 - Re-represents data using the new basis vectors (change of basis)
 - Throws away "uninteresting" directions without much loss of info. This gives a new lower-dimensional representation of the original data
- PCA uses "amount of data variance captured" to define important directions

Machine Learning (CS771A)

・ 同 ト ・ ヨ ト ・ ヨ ト

• How does PCA find the "maximum variance" directions (last class)?

ъ.

イロン 不同と 不同と 不同と

- How does PCA *find* the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data

ъ.

・ロン ・雪と ・ヨン・

- How does PCA *find* the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data
 - Each eigenvector represents one such direction

ъ.

・ロト ・ 一 ・ ・ ヨト ・ ヨト

- How does PCA find the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data
 - Each eigenvector represents one such direction
 - First (top) eigenvector is the direction that captures the largest variance

ъ.

・ロト ・ 一 ・ ・ ヨト ・ ヨト

- How does PCA find the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data
 - Each eigenvector represents one such direction
 - First (top) eigenvector is the direction that captures the largest variance
 - Each subsequent eigenvector is the next best as per this criterion

ヘロト 人間ト 人口ト 人口ト

- How does PCA find the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data
 - Each eigenvector represents one such direction
 - First (top) eigenvector is the direction that captures the largest variance
 - Each subsequent eigenvector is the next best as per this criterion
- Steps in Principal Component Analysis
 - Compute the covariance matrix **S** using the centered data as

$$\mathbf{S} = rac{1}{N} \mathbf{X} \mathbf{X}^ op$$
 (note: **X** assumed $D imes N$ here)

-

- How does PCA find the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data
 - Each eigenvector represents one such direction
 - First (top) eigenvector is the direction that captures the largest variance
 - Each subsequent eigenvector is the next best as per this criterion
- Steps in Principal Component Analysis
 - Compute the covariance matrix **S** using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^{ op}$$
 (note: **X** assumed $D \times N$ here)

• Do an eigen-decomposition of S. This will give D eigenvectors.

-

- How does PCA find the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data
 - Each eigenvector represents one such direction
 - First (top) eigenvector is the direction that captures the largest variance
 - Each subsequent eigenvector is the next best as per this criterion
- Steps in Principal Component Analysis
 - Compute the covariance matrix **S** using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^{ op}$$
 (note: **X** assumed $D \times N$ here)

- Do an eigen-decomposition of **S**. This will give D eigenvectors.
- Take top K leading eigenvectors $\{\boldsymbol{u}_k\}_{k=1}^{K}$ with eigenvalues $\{\lambda_k\}_{k=1}^{K}$

(ロ) (雪) (ヨ) (ヨ) (コ)

- How does PCA find the "maximum variance" directions (last class)?
 - Using eigen-decomposition of the covariance matrix of data
 - Each eigenvector represents one such direction
 - First (top) eigenvector is the direction that captures the largest variance
 - Each subsequent eigenvector is the next best as per this criterion
- Steps in Principal Component Analysis
 - Compute the covariance matrix **S** using the centered data as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^{ op}$$
 (note: **X** assumed $D \times N$ here)

- Do an eigen-decomposition of **S**. This will give D eigenvectors.
- Take top K leading eigenvectors $\{\boldsymbol{u}_k\}_{k=1}^K$ with eigenvalues $\{\lambda_k\}_{k=1}^K$
- $\mathbf{U} = [u_1 \dots u_K]$ is $D \times K$ matrix (each column is a projection direction)

PCA as Linear Projection

• Can use ${f U}$ to linearly project each ${m x}_n \in \mathbb{R}^D$ to a K-dim subspace as



メポト イヨト イヨト

PCA as Linear Projection

• Can use ${f U}$ to linearly project each ${m x}_n \in \mathbb{R}^D$ to a K-dim subspace as



• $\boldsymbol{z}_n \in \mathbb{R}^K$ is also called low-dimensional "embedding" of $\boldsymbol{x}_n \in \mathbb{R}^D$

Machine Learning (CS771A)

・ 同 ト ・ ヨ ト ・ ヨ ト

PCA as Linear Projection

• $\mathbf{Z} = [\mathbf{z}_1 \ \mathbf{z}_2 \dots \ \mathbf{z}_n]$ is the $K \times N$ matrix of embeddings of all the N examples



 $\bullet~{\bf Z}$ can also be thought of as a new, compact feature representation of ${\bf X}$

・ロン ・回と ・ヨン

PCA based Embeddings of Handwritten Digits

• Shown below are 2-dim embeddings of PCA on a set of handwritten digits (each digit image was originally 8×8 , i.e., 64 dimensional)



• PCA to K-dims is also akin to saying $\boldsymbol{x}_n \approx \sum_{k=1}^{K} z_{nk} \boldsymbol{u}_k$. Thus

 $\mathbf{X} \approx \mathbf{UZ}$ (matrix factorization)



3

・ロン ・雪と ・ヨン・

• PCA to K-dims is also akin to saying $\boldsymbol{x}_n \approx \sum_{k=1}^{K} z_{nk} \boldsymbol{u}_k$. Thus

 $\mathbf{X} \approx \mathbf{UZ}$ (matrix factorization)



• Example: Each face in a collection can be represented as a combination of a small no of "eigenfaces" ("template" faces)

ヘロン 人間 とくほ とくほ とうほう

• PCA to K-dims is also akin to saying $\boldsymbol{x}_n \approx \sum_{k=1}^{K} z_{nk} \boldsymbol{u}_k$. Thus

 $\mathbf{X} \approx \mathbf{U}\mathbf{Z}$ (matrix factorization)



• Example: Each face in a collection can be represented as a combination of a small no of "eigenfaces" ("template" faces)

$$\begin{pmatrix} \mathbf{X} \ (\mathbf{D}\mathbf{X}\mathbf{N}) & \mathbf{U} \ (\mathbf{D}\mathbf{X}\mathbf{K}) & \mathbf{Z} \ (\mathbf{K}\mathbf{X}\mathbf{N}) \\ \begin{pmatrix} \mathbf{Q} & \cdots & \mathbf{Q} \end{pmatrix} & \geq \begin{pmatrix} \mathbf{Q} \ (\mathbf{Q} \ \mathbf{Q} \ \mathbf{Q} \ \mathbf{Q} \ \mathbf{Q} & \mathbf{Q} \end{pmatrix} \begin{pmatrix} \mathbf{Z} \ \mathbf{Q} \ \mathbf{Z} \ \mathbf{Q} & \mathbf{Q} \end{pmatrix}$$

Can thus approximately reconstruct the matrix X using UZ: Do PCA on the N × D data matrix X, keep U (D × K) and Z (K × N) and throw away X

ヘロン 不良 とくほど 不良 とうほう

• PCA to K-dims is also akin to saying $\boldsymbol{x}_n \approx \sum_{k=1}^{K} z_{nk} \boldsymbol{u}_k$. Thus

 $\mathbf{X} \approx \mathbf{U}\mathbf{Z}$ (matrix factorization)



• Example: Each face in a collection can be represented as a combination of a small no of "eigenfaces" ("template" faces)

$$\begin{pmatrix} \mathbf{X} (\mathbf{D} \mathbf{X} \mathbf{N}) & \mathbf{U} (\mathbf{D} \mathbf{X} \mathbf{K}) \\ (\mathbf{Q} \dots \mathbf{Q}) \end{pmatrix} \approx \begin{pmatrix} \mathbf{U} (\mathbf{D} \mathbf{X} \mathbf{K}) & \mathbf{Z} (\mathbf{K} \mathbf{X} \mathbf{N}) \\ \mathbf{Z}_1 \dots \mathbf{Z}_n \end{pmatrix}$$

- Can thus approximately reconstruct the matrix X using UZ: Do PCA on the N × D data matrix X, keep U (D × K) and Z (K × N) and throw away X
 - Substantial storage saving if $K \ll D$

◆□▶ ◆□▶ ◆∃▶ ◆∃▶ = うへの

 \bullet Consider doing PCA on a (words \times documents) matrix ${\boldsymbol X}$

-

ヘロト 人間ト 人口ト 人口ト

- $\bullet\,$ Consider doing PCA on a (words $\times\,$ documents) matrix ${\boldsymbol X}$
- Each entry X_{dn} in **X** is the frequency of word d in document n

-

ヘロト 人間ト 人口ト 人口ト

- $\bullet\,$ Consider doing PCA on a (words $\times\,$ documents) matrix ${\boldsymbol X}$
- Each entry X_{dn} in **X** is the frequency of word d in document n
- $\bullet\,$ PCA on ${\bf X}$ will give us eigenvectors that correspond to "topics" or concepts

・ロン ・雪と ・ヨン・

- $\bullet\,$ Consider doing PCA on a (words $\times\,$ documents) matrix ${\boldsymbol X}$
- Each entry X_{dn} in **X** is the frequency of word d in document n
- $\bullet\,$ PCA on ${\bf X}$ will give us eigenvectors that correspond to "topics" or concepts



• Each document is like a weighted combination of these topics

- $\bullet\,$ Consider doing PCA on a (words $\times\,$ documents) matrix ${\boldsymbol X}$
- Each entry X_{dn} in **X** is the frequency of word d in document n
- $\bullet\,$ PCA on ${\bf X}$ will give us eigenvectors that correspond to "topics" or concepts



- Each document is like a weighted combination of these topics
- This is similar to "Latent Semantic Analysis" (LSA), a well-known document dimensionality reduction technique in information retrieval (basically, LSA does SVD on X, which is equivalent to doing PCA on X)

• X is (genes × samples) matrix: Each sample (expression values of a set of D genes) is a weighted combination of K biological "pathways"



글 > - < 글 >

• X is (genes × samples) matrix: Each sample (expression values of a set of D genes) is a weighted combination of K biological "pathways"



• X is (movies × users) matrix: Each user (represented by the vector of his/her ratings of D movies) is a weighted combination of K genres



Beyond Linear Projections..

• Consider the swiss-roll dataset (points lying close to a manifold)



• Linear projection methods (e.g., PCA) can't capture intrinsic nonlinearities

Machine Learning (CS771A)

(日) (同) (日) (日)

Nonlinear Dimensionality Reduction

- Given: Low-dim. surface embedded nonlinearly in high-dim. space
 - Such a structure is called a Manifold



• Goal: Recover the low-dimensional surface



(日) (同) (日) (日)

Nonlinear Dimensionality Reduction

• We want to a learn nonlinear low-dim projection



イロン 不同と イヨン イヨン

Nonlinear Dimensionality Reduction

• We want to a learn nonlinear low-dim projection



• Usually two ways of doing this

(日) (同) (日) (日)
Nonlinear Dimensionality Reduction

• We want to a learn nonlinear low-dim projection



- Usually two ways of doing this
 - Nonlinearize a linear dimensionality reduction method. E.g.:
 - Kernel PCA (nonlinear PCA)

(日) (同) (日) (日)

Nonlinear Dimensionality Reduction

• We want to a learn nonlinear low-dim projection



- Usually two ways of doing this
 - Nonlinearize a linear dimensionality reduction method. E.g.:
 - Kernel PCA (nonlinear PCA)
 - Using manifold based methods that intrinsically preserve nonlinear geometry
 - Locally Linear Embedding (LLE)
 - Isomap
 - Maximum Variance Unfolding
 - Laplacian Eigenmaps
 - And others (tSNE, Hessian LLE, etc.)

イロト イポト イヨト イヨト

Recall PCA: Given N observations {x₁,..., x_N}, ∀x_n ∈ ℝ^D, we define the D × D covariance matrix (assuming centered data ∑_nx_n = 0)

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^{T}$$

ъ.

イロン 不同と 不同と 不同と

Recall PCA: Given N observations {x₁,..., x_N}, ∀x_n ∈ ℝ^D, we define the D × D covariance matrix (assuming centered data ∑_nx_n = 0)

$$\mathbf{S} = rac{1}{N}\sum_{n=1}^N oldsymbol{x}_noldsymbol{x}_n^{ op}$$

• PCA computes eigenvectors \boldsymbol{u}_i which satisfy $\boldsymbol{S}\boldsymbol{u}_i = \lambda_i \boldsymbol{u}_i \ \forall i = 1, \dots, D$

-

ヘロン 人間 とくほ とくほう

Recall PCA: Given N observations {x₁,..., x_N}, ∀x_n ∈ ℝ^D, we define the D × D covariance matrix (assuming centered data ∑_nx_n = 0)

$$\mathbf{S} = rac{1}{N}\sum_{n=1}^N oldsymbol{x}_noldsymbol{x}_n^{ op}$$

- PCA computes eigenvectors \boldsymbol{u}_i which satisfy $\boldsymbol{S}\boldsymbol{u}_i = \lambda_i \boldsymbol{u}_i \; \forall i = 1, \dots, D$
- Let's assume a kernel k with associated M dimensional nonlinear map ϕ

イロン 不良 とくほう 不良 とうせい

Recall PCA: Given N observations {x₁,..., x_N}, ∀x_n ∈ ℝ^D, we define the D × D covariance matrix (assuming centered data ∑_nx_n = 0)

$$\mathbf{S} = rac{1}{N}\sum_{n=1}^N oldsymbol{x}_noldsymbol{x}_n^{ op}$$

- PCA computes eigenvectors \boldsymbol{u}_i which satisfy $\boldsymbol{S}\boldsymbol{u}_i = \lambda_i \boldsymbol{u}_i \; \forall i = 1, \dots, D$
- Let's assume a kernel k with associated M dimensional nonlinear map ϕ
- $M \times M$ covariance matrix in this space (assume centered data $\sum_{n} \phi(\mathbf{x}_{n}) = \mathbf{0}$)

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^{N} \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^{\top}$$

イロン 不良 とくほう イロン しゅう

Recall PCA: Given N observations {x₁,..., x_N}, ∀x_n ∈ ℝ^D, we define the D × D covariance matrix (assuming centered data ∑_nx_n = 0)

$$\mathbf{S} = rac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \mathbf{x}_n^{\top}$$

- PCA computes eigenvectors \boldsymbol{u}_i which satisfy $\boldsymbol{S}\boldsymbol{u}_i = \lambda_i \boldsymbol{u}_i \; \forall i = 1, \dots, D$
- Let's assume a kernel k with associated M dimensional nonlinear map ϕ
- $M \times M$ covariance matrix in this space (assume centered data $\sum_{n} \phi(\mathbf{x}_{n}) = \mathbf{0}$)

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^{N} \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^{\top}$$

• Kernel PCA: Compute eigenvectors \boldsymbol{v}_i satisfying: $C\boldsymbol{v}_i = \lambda_i \boldsymbol{v}_i \ \forall i = 1, \dots, M$

ヘロン 不良 とくほど 不良 とうほう

Recall PCA: Given N observations {x₁,..., x_N}, ∀x_n ∈ ℝ^D, we define the D × D covariance matrix (assuming centered data ∑_nx_n = 0)

$$\mathbf{S} = rac{1}{N}\sum_{n=1}^N oldsymbol{x}_noldsymbol{x}_n^{ op}$$

- PCA computes eigenvectors \boldsymbol{u}_i which satisfy $\boldsymbol{S}\boldsymbol{u}_i = \lambda_i \boldsymbol{u}_i \; \forall i = 1, \dots, D$
- Let's assume a kernel k with associated M dimensional nonlinear map ϕ
- $M \times M$ covariance matrix in this space (assume centered data $\sum_{n} \phi(\mathbf{x}_{n}) = \mathbf{0}$)

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^{N} \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^{\top}$$

- Kernel PCA: Compute eigenvectors \boldsymbol{v}_i satisfying: $\boldsymbol{C}\boldsymbol{v}_i = \lambda_i \boldsymbol{v}_i \ \forall i = 1, \dots, M$
- We would like to do this without having to compute **C** or $\phi(\mathbf{x}_n)$'s

Machine Learning (CS771A)

(日) (四) (日) (日) (日)



Right figure: After mapping the data via ϕ , data is now close to a linear subspace

< ロ > < 同 > < 回 > < 回 >

• Goal: Compute eigenvectors \boldsymbol{v}_i , i.e., $\boldsymbol{C}\boldsymbol{v}_i = \lambda_i \boldsymbol{v}_i$, each \boldsymbol{v}_i is M dimensional

イロン 不得 とくほど 不良 とうほう

- Goal: Compute eigenvectors v_i , i.e., $Cv_i = \lambda_i v_i$, each v_i is M dimensional
- Plugging in the expression for **C**, we have

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\lambda_{i}\boldsymbol{v}$$

(日) (同) (王) (王) (王)

- Goal: Compute eigenvectors v_i , i.e., $Cv_i = \lambda_i v_i$, each v_i is M dimensional
- Plugging in the expression for **C**, we have

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\lambda_{i}\boldsymbol{v}_{i}$$

• Denoting $\alpha_{in} = \frac{1}{\lambda_i N} \phi(\mathbf{x}_n)^\top \mathbf{v}_i$, $\mathbf{v}_i = \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_n)$ (also recall Rep. Thm.)

ヘロン 人間 とくほ とくほ とうほう

- Goal: Compute eigenvectors v_i , i.e., $Cv_i = \lambda_i v_i$, each v_i is M dimensional
- Plugging in the expression for **C**, we have

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\lambda_{i}\boldsymbol{v}_{i}$$

- Denoting $\alpha_{in} = \frac{1}{\lambda_i N} \phi(\mathbf{x}_n)^\top \mathbf{v}_i$, $\mathbf{v}_i = \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_n)$ (also recall Rep. Thm.)
- Thus we can get \boldsymbol{v}_i by finding $\boldsymbol{\alpha}_i = [\alpha_{i1} \ \ldots \ \alpha_{iN}]$

ヘロン 不良 とくほど 不良 とうほう

- Goal: Compute eigenvectors v_i , i.e., $Cv_i = \lambda_i v_i$, each v_i is M dimensional
- Plugging in the expression for **C**, we have

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\lambda_{i}\boldsymbol{v}_{i}$$

• Denoting $\alpha_{in} = \frac{1}{\lambda_i N} \phi(\mathbf{x}_n)^\top \mathbf{v}_i$, $\mathbf{v}_i = \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_n)$ (also recall Rep. Thm.)

- Thus we can get \boldsymbol{v}_i by finding $\boldsymbol{\alpha}_i = [\alpha_{i1} \ \ldots \ \alpha_{iN}]$
- Plugging this back in the eigenvector equation $\mathbf{C}\mathbf{v}_i = \lambda_i \mathbf{v}_i$

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\mathbf{x}_n)\phi(\mathbf{x}_n)^{\top}\sum_{m=1}^{N}\alpha_{im}\phi(\mathbf{x}_m)=\lambda_i\sum_{n=1}^{N}\alpha_{in}\phi(\mathbf{x}_n)$$

・ロン ・日 ・ モン・ モリ・ トロ・

- Goal: Compute eigenvectors v_i , i.e., $Cv_i = \lambda_i v_i$, each v_i is M dimensional
- Plugging in the expression for **C**, we have

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\frac{1}{N}\sum_{n=1}^{N}\phi(\boldsymbol{x}_{n})\phi(\boldsymbol{x}_{n})^{\top}\boldsymbol{v}_{i}=\lambda_{i}\boldsymbol{v}_{i}$$

• Denoting $\alpha_{in} = \frac{1}{\lambda_i N} \phi(\mathbf{x}_n)^\top \mathbf{v}_i$, $\mathbf{v}_i = \sum_{n=1}^N \alpha_{in} \phi(\mathbf{x}_n)$ (also recall Rep. Thm.)

- Thus we can get \boldsymbol{v}_i by finding $\boldsymbol{\alpha}_i = [\alpha_{i1} \ \ldots \ \alpha_{iN}]$
- Plugging this back in the eigenvector equation $\mathbf{C}\mathbf{v}_i = \lambda_i \mathbf{v}_i$

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\mathbf{x}_n)\phi(\mathbf{x}_n)^{\top}\sum_{m=1}^{N}\alpha_{im}\phi(\mathbf{x}_m)=\lambda_i\sum_{n=1}^{N}\alpha_{in}\phi(\mathbf{x}_n)$$

• Pre-multiplying both sides by $\phi(\mathbf{x}_{\ell})^{\top}$ and re-arranging

$$\frac{1}{N}\sum_{n=1}^{N}\phi(\mathbf{x}_{\ell})^{\top}\phi(\mathbf{x}_{n})\sum_{m=1}^{N}\alpha_{im}\phi(\mathbf{x}_{n})^{\top}\phi(\mathbf{x}_{m}) = \lambda_{i}\sum_{n=1}^{N}\alpha_{in}\phi(\mathbf{x}_{\ell})^{\top}\phi(\mathbf{x}_{n})$$

• Using $\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$, we get $\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$

(日本)(四本)(日本)(日本)(日本)

• Using
$$\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$
, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$$

• Define **K** as the $N \times N$ kernel matrix with $K_{nm} = k(\boldsymbol{x}_n, \boldsymbol{x}_m)$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

• Using
$$\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$
, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$$

• Define **K** as the $N \times N$ kernel matrix with $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$

• K is the similarity of two examples x_n and x_m in the ϕ space

3

ヘロト 人間 とくほ とくほ とう

• Using
$$\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$
, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$$

• Define **K** as the $N \times N$ kernel matrix with $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$

- **K** is the similarity of two examples x_n and x_m in the ϕ space
- ϕ is implicitly defined by kernel function k (which can be, e.g., RBF kernel)

イロン 不得 とくほど 不良 とうほう

• Using
$$\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$
, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$$

• Define **K** as the $N \times N$ kernel matrix with $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$

- K is the similarity of two examples x_n and x_m in the ϕ space
- ϕ is implicitly defined by kernel function k (which can be, e.g., RBF kernel)
- Define α_i as the $N \times 1$ vector with elements α_{in}

ヘロン 不良 とくほど 不良 とうほう

• Using
$$\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$
, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$$

• Define **K** as the $N \times N$ kernel matrix with $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$

- K is the similarity of two examples x_n and x_m in the ϕ space
- ϕ is implicitly defined by kernel function k (which can be, e.g., RBF kernel)
- Define α_i as the $N \times 1$ vector with elements α_{in}
- Using \boldsymbol{K} and $\boldsymbol{\alpha}_i$, the eigenvector equation becomes:

$$\mathbf{K}^2 \boldsymbol{\alpha}_i = \lambda_i N \mathbf{K} \boldsymbol{\alpha}_i \quad \Rightarrow \quad \mathbf{K} \boldsymbol{\alpha}_i = \lambda_i N \boldsymbol{\alpha}_i$$

-

・ロト ・ 一 ・ ・ ヨト ・ ヨト

• Using
$$\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$
, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$$

• Define **K** as the $N \times N$ kernel matrix with $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$

- K is the similarity of two examples x_n and x_m in the ϕ space
- ϕ is implicitly defined by kernel function k (which can be, e.g., RBF kernel)
- Define α_i as the $N \times 1$ vector with elements α_{in}
- Using \boldsymbol{K} and $\boldsymbol{\alpha}_i$, the eigenvector equation becomes:

$$\boldsymbol{K}^{2}\boldsymbol{\alpha}_{i}=\lambda_{i}\boldsymbol{N}\boldsymbol{K}\boldsymbol{\alpha}_{i}$$
 \Rightarrow $\boldsymbol{K}\boldsymbol{\alpha}_{i}=\lambda_{i}\boldsymbol{N}\boldsymbol{\alpha}_{i}$

• Thus α_i is an eigenvector of the N imes N kernel matrix **K**

イロン 不良 とくほう イロン しゅう

• Using
$$\phi(\mathbf{x}_n)^{\top} \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$$
, we get

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_{\ell}, \mathbf{x}_n) \sum_{m=1}^N \alpha_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N \alpha_{in} k(\mathbf{x}_{\ell}, \mathbf{x}_n)$$

• Define **K** as the $N \times N$ kernel matrix with $K_{nm} = k(\mathbf{x}_n, \mathbf{x}_m)$

- K is the similarity of two examples x_n and x_m in the ϕ space
- ϕ is implicitly defined by kernel function k (which can be, e.g., RBF kernel)
- Define α_i as the $N \times 1$ vector with elements α_{in}
- Using \boldsymbol{K} and $\boldsymbol{\alpha}_i$, the eigenvector equation becomes:

$$\mathbf{K}^2 \boldsymbol{\alpha}_i = \lambda_i N \mathbf{K} \boldsymbol{\alpha}_i \quad \Rightarrow \quad \mathbf{K} \boldsymbol{\alpha}_i = \lambda_i N \boldsymbol{\alpha}_i$$

- Thus α_i is an eigenvector of the N imes N kernel matrix **K**
- Note: Since $\mathbf{v}_i^{\top} \mathbf{v}_i = 1$ and $\mathbf{v}_i = \sum_{n=1}^{N} \alpha_{in} \phi(\mathbf{x}_n)$, we have $\alpha_i^{\top} \mathbf{K} \alpha_i = 1$, which means $\alpha_i^{\top} \lambda_i N \alpha_i = 1$ and $\alpha_i^{\top} \alpha_i = 1/(\lambda_i N)$. Thus the original solution with $\alpha_i^{\top} \alpha_i = 1$ (eigenvec with norm 1) needs to be re-scaled as $\tilde{\alpha}_{in} = \alpha_{in}/\sqrt{\lambda_i N}$

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$\widetilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - \frac{1}{N} \sum_{\ell=1}^{N} \phi(\boldsymbol{x}_\ell)$$

-

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$ilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - rac{1}{N}\sum_{\ell=1}^N \phi(\boldsymbol{x}_\ell)$$

• Each element of the centered kernel matrix

ъ.

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$ilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - \frac{1}{N} \sum_{\ell=1}^N \phi(\boldsymbol{x}_\ell)$$

• Each element of the centered kernel matrix

$$ilde{K}_{nm} = ilde{\phi}(\mathbf{x}_n)^{\top} ilde{\phi}(\mathbf{x}_m)$$

ъ.

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$\widetilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - \frac{1}{N} \sum_{\ell=1}^{N} \phi(\boldsymbol{x}_\ell)$$

• Each element of the centered kernel matrix

$$\begin{split} \tilde{\mathcal{K}}_{nm} &= \quad \tilde{\phi}(\boldsymbol{x}_n)^\top \tilde{\phi}(\boldsymbol{x}_m) \\ &= \quad \phi(\boldsymbol{x}_n)^\top \phi(\boldsymbol{x}_m) - \frac{1}{N} \sum_{\ell=1}^N \phi(\boldsymbol{x}_n)^\top \phi(\boldsymbol{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N \phi(\boldsymbol{x}_\ell)^\top \phi(\boldsymbol{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N \phi(\boldsymbol{x}_j)^\top \phi(\boldsymbol{x}_\ell) \end{split}$$

ъ.

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$\widetilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - \frac{1}{N} \sum_{\ell=1}^{N} \phi(\boldsymbol{x}_\ell)$$

• Each element of the centered kernel matrix

$$\begin{split} \tilde{\mathcal{K}}_{nm} &= \tilde{\phi}(\mathbf{x}_n)^\top \tilde{\phi}(\mathbf{x}_m) \\ &= \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_\ell)^\top \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_\ell) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_n, \mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_\ell) \end{split}$$

ъ.

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$\widetilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - \frac{1}{N} \sum_{\ell=1}^N \phi(\boldsymbol{x}_\ell)$$

• Each element of the centered kernel matrix

$$\begin{split} \tilde{\mathcal{K}}_{nm} &= \tilde{\phi}(\mathbf{x}_n)^\top \tilde{\phi}(\mathbf{x}_m) \\ &= \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_\ell)^\top \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_\ell) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_n, \mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_\ell) \end{split}$$

 \bullet In matrix notation, the centered $\tilde{\textbf{\textit{K}}}=\textbf{\textit{K}}-\textbf{1}_{N}\textbf{\textit{K}}-\textbf{\textit{K}}\textbf{1}_{N}+\textbf{1}_{N}\textbf{\textit{K}}\textbf{1}_{N}$

-

ヘロン 人間 とくほ とくほう

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$\widetilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - \frac{1}{N} \sum_{\ell=1}^N \phi(\boldsymbol{x}_\ell)$$

• Each element of the centered kernel matrix

$$\begin{split} \tilde{K}_{nm} &= \tilde{\phi}(\mathbf{x}_n)^\top \tilde{\phi}(\mathbf{x}_m) \\ &= \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_\ell)^\top \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_\ell) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_n, \mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_\ell) \end{split}$$

 \bullet In matrix notation, the centered $\tilde{\textbf{\textit{K}}}=\textbf{\textit{K}}-\textbf{1}_{\textit{N}}\textbf{\textit{K}}-\textbf{\textit{K}}\textbf{1}_{\textit{N}}+\textbf{1}_{\textit{N}}\textbf{\textit{K}}\textbf{1}_{\textit{N}}$

• $\mathbf{1}_N$ is the $N \times N$ matrix with every element = 1/N

Machine Learning (CS771A)

- In PCA, we centered the data before computing the covariance matrix
- For kernel PCA, we need to do the same

$$\tilde{\phi}(\boldsymbol{x}_n) = \phi(\boldsymbol{x}_n) - \frac{1}{N} \sum_{\ell=1}^{N} \phi(\boldsymbol{x}_\ell)$$

• Each element of the centered kernel matrix

$$\begin{split} \tilde{K}_{nm} &= \tilde{\phi}(\mathbf{x}_n)^\top \tilde{\phi}(\mathbf{x}_m) \\ &= \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_n)^\top \phi(\mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N \phi(\mathbf{x}_\ell)^\top \phi(\mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_\ell) \\ &= k(\mathbf{x}_n, \mathbf{x}_m) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_n, \mathbf{x}_\ell) - \frac{1}{N} \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_m) + \frac{1}{N^2} \sum_{j=1}^N \sum_{\ell=1}^N k(\mathbf{x}_\ell, \mathbf{x}_\ell) \end{split}$$

- \bullet In matrix notation, the centered $\tilde{\textbf{\textit{K}}}=\textbf{\textit{K}}-\textbf{1}_{\textit{N}}\textbf{\textit{K}}-\textbf{\textit{K}}\textbf{1}_{\textit{N}}+\textbf{1}_{\textit{N}}\textbf{\textit{K}}\textbf{1}_{\textit{N}}$
- $\mathbf{1}_N$ is the N imes N matrix with every element = 1/N
- Eigen-decomposition is then done for the centered kernel matrix $ilde{m{K}}$

Machine Learning (CS771A)

- Suppose $\{\alpha_1, \ldots, \alpha_L\}$ are the top L eigenvectors of kernel matrix $ilde{m{\kappa}}$
- The L-dimensional KPCA projection $\boldsymbol{z}_m = [z_{m1}, \dots, z_{mL}]$ of a point \boldsymbol{x}_m :

$$z_{m\ell} = \phi(\boldsymbol{x}_m)^\top \boldsymbol{v}_\ell \qquad \forall \ell = 1, \dots, L$$

イロン 不良 とくほう イロン しゅう

- Suppose $\{ lpha_1, \dots, lpha_L \}$ are the top L eigenvectors of kernel matrix $ilde{m{\kappa}}$
- The *L*-dimensional KPCA projection $\boldsymbol{z}_m = [z_{m1}, \dots, z_{mL}]$ of a point \boldsymbol{x}_m :

$$z_{m\ell} = \phi(\boldsymbol{x}_m)^\top \boldsymbol{v}_\ell \qquad \forall \ell = 1, \dots, L$$

• Using the definition of \boldsymbol{v}_{ℓ} , i.e., $\boldsymbol{v}_{\ell} = \sum_{n=1}^{N} \alpha_{\ell n} \phi(\boldsymbol{x}_n)$, we have

$$z_{m\ell} = \phi(\boldsymbol{x}_m)^{\top} \boldsymbol{v}_{\ell} = \sum_{n=1}^{N} \alpha_{\ell n} k(\boldsymbol{x}_n, \boldsymbol{x}_m)$$

イロン 不良 とくほう イロン しゅう

- Suppose $\{ lpha_1, \ldots, lpha_L \}$ are the top L eigenvectors of kernel matrix $ilde{m{\kappa}}$
- The *L*-dimensional KPCA projection $\boldsymbol{z}_m = [z_{m1}, \dots, z_{mL}]$ of a point \boldsymbol{x}_m :

$$z_{m\ell} = \phi(\boldsymbol{x}_m)^\top \boldsymbol{v}_\ell \qquad \forall \ell = 1, \dots, L$$

• Using the definition of \boldsymbol{v}_{ℓ} , i.e., $\boldsymbol{v}_{\ell} = \sum_{n=1}^{N} \alpha_{\ell n} \phi(\boldsymbol{x}_n)$, we have

$$z_{m\ell} = \phi(\boldsymbol{x}_m)^{\top} \boldsymbol{v}_{\ell} = \sum_{n=1}^{N} \alpha_{\ell n} k(\boldsymbol{x}_n, \boldsymbol{x}_m)$$

• Note: Cost of computing the embeddings scales in N

- Suppose $\{ lpha_1, \dots, lpha_L \}$ are the top L eigenvectors of kernel matrix $ilde{m{\kappa}}$
- The *L*-dimensional KPCA projection $\boldsymbol{z}_m = [z_{m1}, \dots, z_{mL}]$ of a point \boldsymbol{x}_m :

$$z_{m\ell} = \phi(\boldsymbol{x}_m)^\top \boldsymbol{v}_\ell \qquad \forall \ell = 1, \dots, L$$

• Using the definition of \boldsymbol{v}_{ℓ} , i.e., $\boldsymbol{v}_{\ell} = \sum_{n=1}^{N} \alpha_{\ell n} \phi(\boldsymbol{x}_n)$, we have

$$z_{m\ell} = \phi(\boldsymbol{x}_m)^{\top} \boldsymbol{v}_{\ell} = \sum_{n=1}^{N} \alpha_{\ell n} k(\boldsymbol{x}_n, \boldsymbol{x}_m)$$

- Note: Cost of computing the embeddings scales in *N*
- Note: For linear kernel, KPCA reduces to PCA (but more efficient if N < D)

Kernel PCA: Summary of the algorithm

• Construct the $N \times N$ kernel matrix **K**

ъ.

イロン 不同と イヨン イヨン
• Construct the $N \times N$ kernel matrix **K**

• Center K as follows $\tilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N$, where $\mathbf{1}_N$ is an $N \times N$ matrix of all 1/N

ъ.

- Construct the $N \times N$ kernel matrix ${\bf K}$
- Center K as follows $\tilde{K} = K \mathbf{1}_N K K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N$, where $\mathbf{1}_N$ is an $N \times N$ matrix of all 1/N
- Do eigen-decomposition of \tilde{K} and find top L eigenvecs $\alpha_1, \alpha_2, \ldots, \alpha_L$ with eigenvals $\lambda_1, \lambda_2, \ldots, \lambda_L$

イロン 不良 とくほう 不良 とうせい

• Construct the $N \times N$ kernel matrix ${\bf K}$

• Center K as follows $\tilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N$, where $\mathbf{1}_N$ is an $N \times N$ matrix of all 1/N

- Do eigen-decomposition of \tilde{K} and find top L eigenvecs $\alpha_1, \alpha_2, \ldots, \alpha_L$ with eigenvals $\lambda_1, \lambda_2, \ldots, \lambda_L$
- Re-scale each eigenvector as $\tilde{m{lpha}}_i = m{lpha}_i / \sqrt{\lambda_i N}$, $\forall i = 1, \dots, L$

• Construct the $N \times N$ kernel matrix ${\bf K}$

• Center K as follows $\tilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N$, where $\mathbf{1}_N$ is an $N \times N$ matrix of all 1/N

- Do eigen-decomposition of \tilde{K} and find top L eigenvecs $\alpha_1, \alpha_2, \ldots, \alpha_L$ with eigenvals $\lambda_1, \lambda_2, \ldots, \lambda_L$
- Re-scale each eigenvector as $ilde{lpha}_i = {lpha}_i / \sqrt{\lambda_i N}$, $\forall i=1,\ldots,L$
- Finally, compute embedding $\boldsymbol{z}_m \in \mathbb{R}^L$ of any point \boldsymbol{x}_m as

$$z_{m\ell} = \sum_{n=1}^{N} \tilde{lpha}_{\ell n} k(\boldsymbol{x}_n, \boldsymbol{x}_m)$$

• Construct the $N \times N$ kernel matrix ${\bf K}$

• Center K as follows $\tilde{K} = K - \mathbf{1}_N K - K \mathbf{1}_N + \mathbf{1}_N K \mathbf{1}_N$, where $\mathbf{1}_N$ is an $N \times N$ matrix of all 1/N

- Do eigen-decomposition of \tilde{K} and find top L eigenvecs $\alpha_1, \alpha_2, \ldots, \alpha_L$ with eigenvals $\lambda_1, \lambda_2, \ldots, \lambda_L$
- Re-scale each eigenvector as $ilde{m{lpha}}_i = {m{lpha}}_i / \sqrt{\lambda_i N}$, $orall i = 1, \dots, L$
- Finally, compute embedding $\boldsymbol{z}_m \in \mathbb{R}^L$ of any point \boldsymbol{x}_m as

$$z_{m\ell} = \sum_{n=1}^{N} \tilde{lpha}_{\ell n} k(oldsymbol{x}_n, oldsymbol{x}_m)$$

Note: For compactness, the L × N matrix of all L eigenvectors (each is N dimensional) can be written as α̃ = [α̃₁ α̃₂ ... α̃_L]^T. Thus we can also write embedding z_m ∈ ℝ^L as z_m = α̃k_m where k_m is N × 1 vector of kernelized similarities of x_m with all the N training data points

うしつ 山田 シスト・エー・ 日 うのつ

Kernel PCA: An Example



Note that even if we throw away the 2nd PC, we get a good 1-D embedding