

Learning from Data

Lecture 10: Principal Component Analysis

Yang Li yangli@sz.tsinghua.edu.cn

TBSI

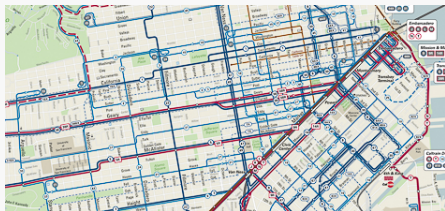
December 2, 2022

Unsupervised Learning (Part II): PCA

- ▶ Motivation
- ▶ Linear PCA
- ▶ Kernel PCA

Motivation of PCA

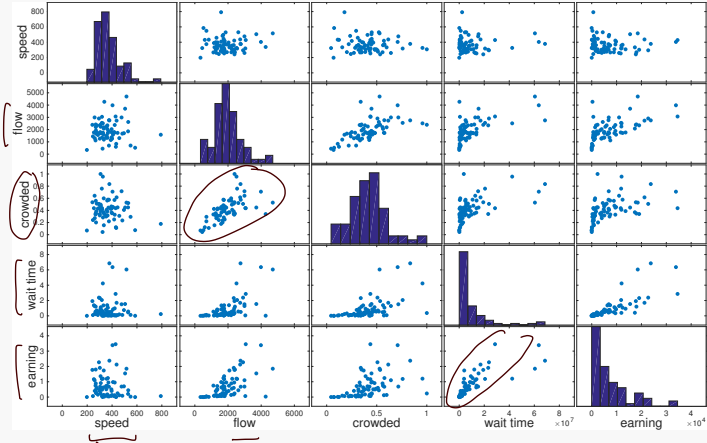
Example: Analyzing San Francisco public transit route efficiency



features	notes
speed	average speed
flow	# boarding passengers per hour
crowded	% passenger capacity reached
wait time	average waiting time at bus stop
earning	net operation revenue
:	:

Motivation of PCA

Input features contain a lot of redundancy

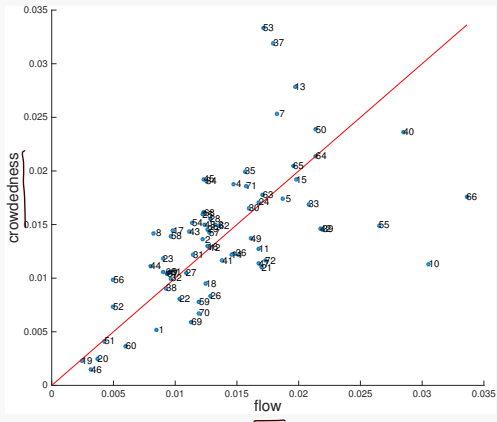


Scatter plot matrix reveals pairwise correlations among 5 major features

Motivation of PCA

Example of linearly dependent features

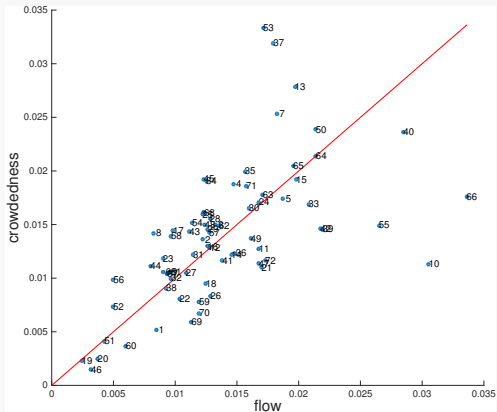
- ▶ Flow: average # boarding passengers per hour
- ▶ Crowdedness: $\frac{\text{average \# passengers on train}}{\text{train capacity}}$



Motivation of PCA

Example of linearly dependent features

- ▶ Flow: average # boarding passengers per hour
- ▶ Crowdedness: $\frac{\text{average \# passengers on train}}{\text{train capacity}}$



How can we automatically detect and remove this redundancy?

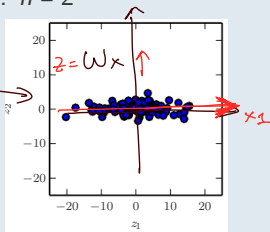
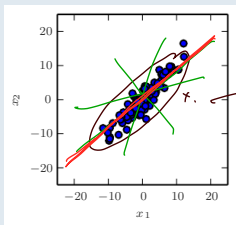
- ▶ geometric approach ← *start here!*
- ▶ diagonalize covariance matrix approach

How to remove feature redundancy?

Given $\{x^{(1)}, \dots, x^{(m)}\}$, $x^{(i)} \in \mathbb{R}^n$.

- ▶ Find a linear, orthogonal transformation $W : \mathbb{R}^n \rightarrow \mathbb{R}^k$ of the input data
- ▶ W aligns the direction of maximum variance with the axes of the new space.

Example: $n = 2$



features x_1 and x_2 are strongly correlated

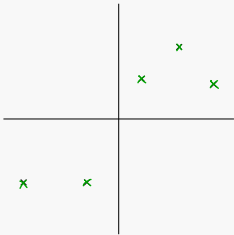
variations in $z = x^T W$ is mostly along the x -axis. x can be represented in 1D! z .

Direction of Maximum Variance

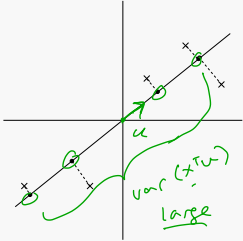
- ▶ Suppose $\mu = \text{mean}(x) = 0$, $\sigma_j = \text{var}(x_j) = 1$ (variance of jth feature)

Direction of Maximum Variance

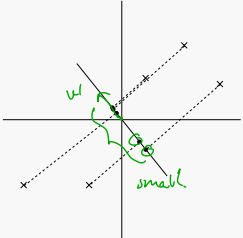
- ▶ Suppose $\mu = \text{mean}(x) = 0$, $\sigma_j = \text{var}(x_j) = 1$ (variance of j th feature)
- ▶ Find major axis of variation unit vector u :



input observations



projections on u have large variance



projections on u have small variance

u maximizes the variance of the projections

Principal Component Analysis (PCA)

Pearson, K. (1901), Hotelling, H. (1933) "Analysis of a complex of statistical variables into principal components". Journal of Educational Psychology.

PCA goals

- ▶ Find principal components u_1, \dots, u_n that are mutually orthogonal (uncorrelated)
- ▶ Most of the variation in x will be accounted for by k principal components where $k \ll n$.

Principal Component Analysis (PCA)

Pearson, K. (1901), Hotelling, H. (1933) "Analysis of a complex of statistical variables into principal components". Journal of Educational Psychology.

PCA goals

- ▶ Find principal components u_1, \dots, u_n that are mutually orthogonal (uncorrelated)
- ▶ Most of the variation in x will be accounted for by k principal components where $k \ll n$.

Main steps of (full) PCA:

1. Standardize x such that $Mean(x) = 0$, $Var(x_j) = 1$ for all j
2. Find projection of x , $u_1^T x$ with maximum variance u_1 : 1st principal component.
3. For $j = 2, \dots, n$,
Find another projection of x , $u_j^T x$ with maximum variance,
where u_j is orthogonal to u_1, \dots, u_{j-1}

Step 1: Standardize data

Normalize x such that Mean(x) = 0 and Var(x_j) = 1

$$x^{(i)} := x^{(i)} - \mu \leftarrow \text{recenter}$$

$$x_j^{(i)} := x_j^{(i)} / \sigma_j \leftarrow \text{scale by } \text{stdev}(x_j)$$

Check:

$$\begin{aligned} \text{var}\left(\frac{x_j}{\sigma_j}\right) &= \frac{1}{m} \sum_{i=1}^m \left(\frac{x_j^{(i)} - \mu_j}{\sigma_j} \right)^2 = \frac{1}{\sigma_j^2} \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2 \\ &= \frac{1}{\sigma_j^2} \sigma_j^2 = 1 \end{aligned}$$

Handwritten notes: $\delta_j = \text{Var}(x_j)$, δ_j^2

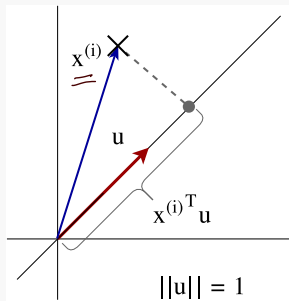
Step 2: Find Projection with Maximum Variance

$$X = \begin{bmatrix} -x^{(1)} \\ -x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix} \left. \vphantom{\begin{bmatrix} -x^{(1)} \\ -x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix}} \right\} m$$

$$\text{COV}(X) = \begin{bmatrix} b_{11}^2 & b_{12}^2 \\ & b_{12}^2 & \dots \\ & & \dots & \dots \\ & & & \dots & b_{nn}^2 \end{bmatrix}$$

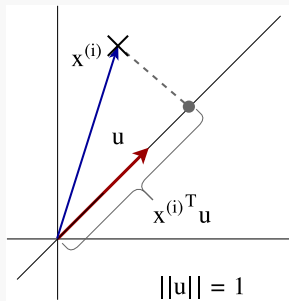
$(n \times n)$

Since $\|u\| = 1$, the length of $x^{(i)}$'s projection on u is $x^{(i)T}u$.



Step 2: Find Projection with Maximum Variance

Since $\|u\| = 1$, the length of $x^{(i)}$'s projection on u is $x^{(i)T}u$.

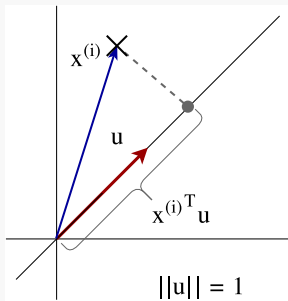


Variance of the projections:

$$\frac{1}{m} \sum_{i=1}^m (x^{(i)T}u - 0)^2 = \frac{1}{m} \sum_{i=1}^m u^T x^{(i)} x^{(i)T} u$$

Step 2: Find Projection with Maximum Variance

Since $\|u\| = 1$, the length of $x^{(i)}$'s projection on u is $x^{(i)T}u$.

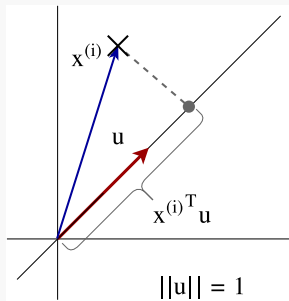


Variance of the projections:

$$\begin{aligned}\frac{1}{m} \sum_{i=1}^m (x^{(i)T}u - 0)^2 &= \frac{1}{m} \sum_{i=1}^m u^T x^{(i)} x^{(i)T} u \\ &= u^T \left(\underbrace{\frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}}_{\text{cov}(x) = \Sigma} \right) u\end{aligned}$$

Step 2: Find Projection with Maximum Variance

Since $\|u\| = 1$, the length of $x^{(i)}$'s projection on u is $x^{(i)T}u$.



Variance of the projections:

$$\begin{aligned}\frac{1}{m} \sum_{i=1}^m (x^{(i)T}u - 0)^2 &= \frac{1}{m} \sum_{i=1}^m u^T x^{(i)} x^{(i)T} u \\ &= u^T \left(\frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \right) u \\ &= u^T \Sigma u\end{aligned}$$

Σ : the sample covariance matrix of $x^{(1)} \dots x^{(m)}$.

1st Principal Component

Find unit vector u_1 that maximizes variance of projections:

$$u_1 = \operatorname{argmax}_{u: \|u\|=1} \underline{u}^T \Sigma u \quad (1)$$

(definition)

u_1 is the **1st principal component** of X

u_1 can be solved using optimization tools, but it has a more efficient solution:

Proposition 1

u_1 is the largest eigenvector of covariance matrix Σ

Proposition 1

u_1 is the largest eigenvector of covariance matrix Σ

Proof.



$$\max_{\|u\|=1} u^T \Sigma u \rightarrow \min_{u^T u = 1} -u^T \Sigma u$$

Proposition 1

u_1 is the largest eigenvector of covariance matrix Σ

Proof. Generalized Lagrange function of Problem **??**:

$$L(u) = -u^T \Sigma u + \beta (u^T u - 1)$$

Proposition 1

u_1 is the largest eigenvector of covariance matrix Σ

Proof. Generalized Lagrange function of Problem ??:

$$L(u) = -u^T \Sigma u + \beta(u^T u - 1)$$

To minimize $L(u)$,

$$\frac{\delta L}{\delta u} = -2\Sigma u + 2\beta u = 0 \implies \underline{\Sigma u = \beta u}$$

Therefore u_1 must be an eigenvector of Σ .

Proposition 1

u_1 is the largest eigenvector of covariance matrix Σ

Proof. Generalized Lagrange function of Problem ??:

$$L(u) = -u^T \Sigma u + \beta(u^T u - 1)$$

To minimize $L(u)$,

$$\frac{\delta L}{\delta u} = -2\Sigma u + 2\beta u = 0 \implies \Sigma u = \beta u$$

Therefore u_1 must be an eigenvector of Σ .

Let $u_1 = v_j$, the eigenvector with the j th largest eigenvalue λ_j ,

$$u_1^T \Sigma u_1 = v_j^T \Sigma v_j = \lambda_j v_j^T v_j = \lambda_j.$$

Hence $\underline{u_1} = \underline{v_1}$, the eigenvector with the largest eigenvalue $\underline{\lambda_1}$. □

Proposition 2

The j th principal component of X , u_j is the j th largest eigenvector of Σ .

Proof.

Proposition 2

The j th principal component of X , u_j is the j th largest eigenvector of Σ .

Proof. Consider the case $j = 2$,

$$u_2 = \operatorname{argmax}_{u: \|u\|=1, u_1^T u=0} u^T \Sigma u \quad (2)$$

Proposition 2

The j th principal component of X , u_j is the j th largest eigenvector of Σ .

Proof. Consider the case $j = 2$,

$$u_2 = \underset{u: \|u\|=1, u_1^T u=0}{\operatorname{argmax}} u^T \Sigma u \quad (2)$$

The Lagrangian function:

$$L(u) = -u^T \Sigma u + \beta_1 (u^T u - 1) + \beta_2 (u_1^T u)$$

Minimizing $L(u)$ yields:

$$\begin{aligned} \beta_2 = 0, \Sigma u = \beta_1 u \\ \frac{\partial L(u)}{\partial u} = -2 \Sigma u + 2\beta_1 u + \beta_2 u_1 = 0. \quad \frac{\partial L(u)}{\partial u} = -2 \Sigma u + 2\beta_1 u = 0. \\ \text{Multiply by } u_1^T: \quad -2 u_1^T \Sigma u + 2\beta_1 u_1^T u + \beta_2 u_1^T u_1 = 0. \quad \Sigma u = \beta_1 u. \\ \Rightarrow \underbrace{\lambda u_1^T u}_{=0} + 2\beta_1 \underbrace{u_1^T u}_{=0} + \beta_2 \underbrace{u_1^T u_1}_{=1} = 0. \\ \left(\frac{\partial L}{\partial \beta_2} = 0 \Rightarrow u_1^T u = 0 \right) \quad \beta_2 = 0. \end{aligned}$$

Proposition 2

The j th principal component of X , u_j is the j th largest eigenvector of Σ .

Proof. Consider the case $j = 2$,

$$u_2 = \operatorname{argmax}_{u: \|u\|=1, u_1^T u=0} u^T \Sigma u \quad (2)$$

The Lagrangian function:

$$L(u) = -u^T \Sigma u + \beta_1 (u^T u - 1) + \beta_2 (u_1^T u)$$

Minimizing $L(u)$ yields:

$$\beta_2 = 0, \Sigma u = \beta_1 u$$

To maximize $\underline{u^T \Sigma u} = \lambda$, u_2 must be the eigenvector with the second largest eigenvalue $\beta_1 = \lambda_2$. The same argument can be generalized to cases $j > 2$. (Use induction to prove for $j = 1 \dots n$) □

Summary

We can solve PCA by solving an eigenvalue problem!

Main steps of (full) PCA:

1. Standardize x such that $Mean(x) = 0$, $Var(x_j) = 1$ for all j
2. Compute $\Sigma = cov(x)$
3. Find principal components u_1, \dots, u_n by eigenvalue decomposition:
 $\Sigma = U\Lambda U^T$. $\leftarrow U$ is an orthogonal basis in \mathbb{R}^n

Next we project data vectors x to this new basis, which spans the **principal component space**.

PCA Projection

- ▶ Projection of sample $\underline{x} \in \mathbb{R}^n$ in the principal component space:

projection
↓

$$\underline{z}^{(i)} = \begin{bmatrix} x^{(i)T} \underline{u}_1 \\ \vdots \\ x^{(i)T} \underline{u}_n \end{bmatrix} \in \mathbb{R}^n$$

- ▶ Projection of sample $x \in \mathbb{R}^n$ in the principal component space:

$$z^{(i)} = \begin{bmatrix} x^{(i)T} u_1 \\ \vdots \\ x^{(i)T} u_n \end{bmatrix} \in \mathbb{R}^n$$

- ▶ Matrix notation:

$$z^{(i)} = \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix}^T x^{(i)} = U^T x^{(i)}, \text{ or } Z = XU$$

- ▶ Projection of sample $x \in \mathbb{R}^n$ in the principal component space:

$$z^{(i)} = \begin{bmatrix} x^{(i)T} u_1 \\ \vdots \\ x^{(i)T} u_n \end{bmatrix} \in \mathbb{R}^n$$

- ▶ Matrix notation:

$$z^{(i)} = \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix}^T x^{(i)} = U^T x^{(i)}, \text{ or } Z = XU$$

- ▶ The truncated transformation $Z_k = XU_k$ keeping only the first k principal components is used for **dimension reduction**.

Properties of PCA

- ▶ The variance of principal component projections are

$$\text{Var}(\underbrace{x^T}_{\text{cov}(x)} \underbrace{u_j}_{\text{cov}(x)}) = \underbrace{u_j^T \Sigma u_j}_{\text{cov}(x)} = \lambda_j \text{ for } j = 1, \dots, n$$

Properties of PCA

- ▶ The variance of principal component projections are

$$\text{Var}(x^T u_j) = \underline{u_j}^T \Sigma u_j = \lambda_j \text{ for } j = 1, \dots, n$$

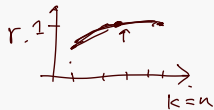
Var

- ▶ % of variance explained by the j th principal component: $\frac{\lambda_j}{\sum_{i=1}^n \lambda_i}$.
i.e. projections are uncorrelated

Properties of PCA

- ▶ The variance of principal component projections are

$$\text{Var}(x^T u_j) = u_j^T \Sigma u_j = \lambda_j \text{ for } j = 1, \dots, n$$



- ▶ % of variance explained by the j th principal component: $\frac{\lambda_j}{\sum_{i=1}^n \lambda_i}$.
i.e. projections are uncorrelated

- r. ▶ % of variance accounted for by retaining the first k principal components ($k \leq n$): $\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^n \lambda_j}$ } var. of first k pc projections.

Another geometric interpretation of PCA is minimizing projection residuals. (see homework!)

Covariance Interpretation of PCA

$$X = \text{cov}(X) = \Sigma$$

PCA removes the “redundancy” (or noise) in input data X:

Let Z = XU be the PCA projected data,

$$\text{cov}(Z) = \frac{1}{m} Z^T Z = \frac{1}{m} \underbrace{(XU)^T}_{\text{cov}(X)} \underbrace{(XU)}_{\text{cov}(X)} = U^T \underbrace{\left(\frac{1}{m} X^T X \right)}_{\Sigma} U = \underbrace{U^T}_{\text{cov}(X)} \underbrace{\Sigma}_{\text{cov}(X)} U$$

Covariance Interpretation of PCA

PCA removes the “redundancy” (or noise) in input data X :

Let $Z = XU$ be the PCA projected data,

$$\text{cov}(Z) = \frac{1}{m} Z^T Z = \frac{1}{m} (XU)^T (XU) = U^T \left(\frac{1}{m} X^T X \right) U = U^T \Sigma U$$

Since Σ is symmetric, it has real eigenvalues. Its eigen decomposition is

$$\Sigma = \underline{U \Lambda U^T}$$

where

$$\underline{U} = \begin{bmatrix} | & & | \\ \underline{u_1} & \dots & \underline{u_n} \\ | & & | \end{bmatrix}, \Lambda = \begin{bmatrix} \underline{\lambda_1} & & \\ & \ddots & \\ & & \underline{\lambda_n} \end{bmatrix}$$

Covariance Interpretation of PCA

PCA removes the “redundancy” (or noise) in input data X :

Let $Z = XU$ be the PCA projected data,

$$\text{cov}(Z) = \frac{1}{m} Z^T Z = \frac{1}{m} (XU)^T (XU) = U^T \left(\frac{1}{m} X^T X \right) U = U^T \Sigma U$$

Since Σ is symmetric, it has real eigenvalues. Its eigen decomposition is

$$\Sigma = U \Lambda U^T$$

where

$$U = \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix}, \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

Then

$$\text{cov}(Z) = U^T \underbrace{(U \Lambda U^T)}_{I} U = \underbrace{\Lambda}_I$$

Covariance Interpretation of PCA

PCA removes the “redundancy” (or noise) in input data X :

Let $Z = XU$ be the PCA projected data,

$$\text{cov}(Z) = \frac{1}{m} Z^T Z = \frac{1}{m} (XU)^T (XU) = U^T \left(\frac{1}{m} X^T X \right) U = U^T \Sigma U$$

Since Σ is symmetric, it has real eigenvalues. Its eigen decomposition is

$$\Sigma = U \Lambda U^T$$

where

$$U = \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix}, \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{bmatrix}$$

Then

$$\text{cov}(Z) = U^T (U \Lambda U^T) U = \Lambda$$

The principal component transformation XU diagonalizes the sample covariance matrix of X

PCA Dimension reduction

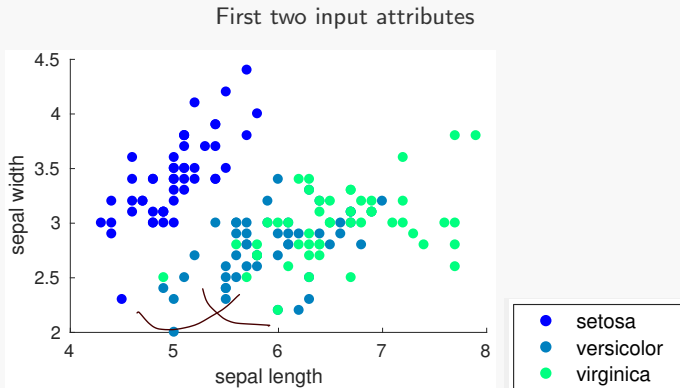
- ▶ Find principal components $\underline{u_1}, \dots, \underline{u_n}$ that are mutually orthogonal (uncorrelated)
- ▶ Most of the variations in x will be accounted for by k principal components where $k \ll n$.

Main steps

1. Standardize x such that $Mean(x) = 0$, $Var(x_j) = 1$ for all j
2. Compute $\underline{\Sigma} = cov(x)$
3. Find principal components u_1, \dots, u_n by eigenvalue decomposition:
 $\underline{\Sigma} = U \Lambda U^T$. $\leftarrow U$ is an orthogonal basis in \mathbb{R}^n
4. Project data on first the k principal components:
 $\underline{z} = [x^T u_1, \dots, x^T u_k]^T$

PCA Example: Iris Dataset

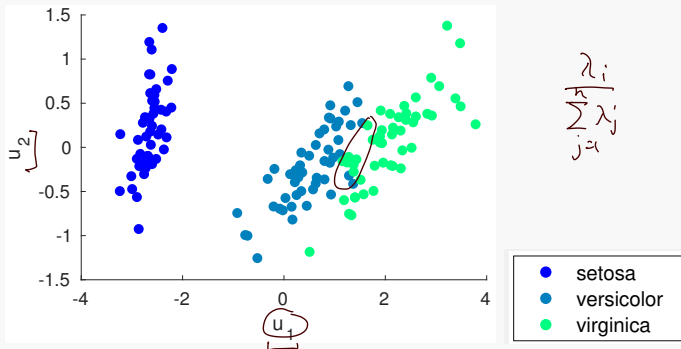
- ▶ 150 samples
- ▶ input feature dimension: 4



PCA Example: Iris Dataset

- ▶ 150 samples
- ▶ input feature dimension: 4

PCA Projection on 2 Principal Components



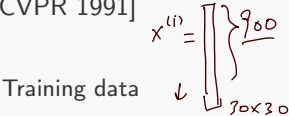
% of variance explained by PC1: 73%,

by PC2: 22%

95%

PCA Example: Eigenfaces

Learning image representations for face recognition using PCA [Turk and Pentland CVPR 1991]

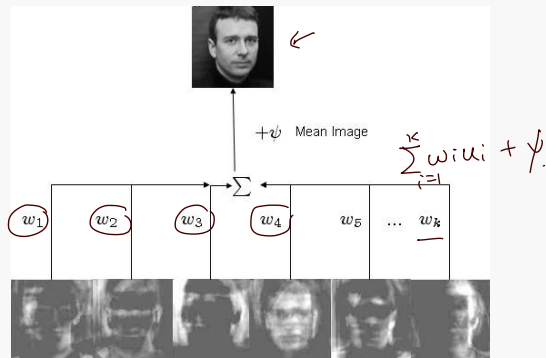


Eigenfaces: k principal components u_1, u_2, u_3



PCA Example: Eigenfaces

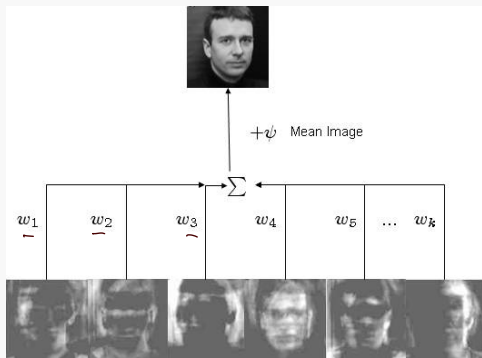
Each face image is a linear combination of the eigenfaces (principal components)



Each image is represented by k weights

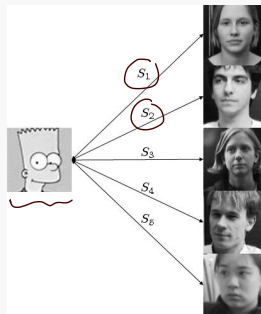
PCA Example: Eigenfaces

Each face image is a linear combination of the **eigenfaces** (principal components)



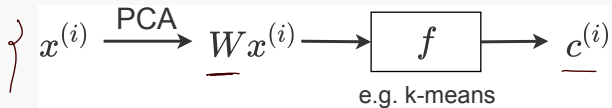
Each image is represented by k weights

Recognize faces by classifying the weight vectors. e.g. k-Nearest Neighbor



Kernel PCA

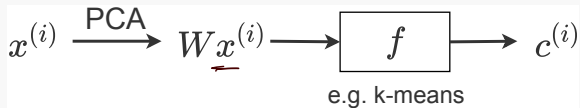
Feature extraction using PCA



Linear PCA assumes data are separable in \mathbb{R}^n

Kernel PCA

Feature extraction using PCA



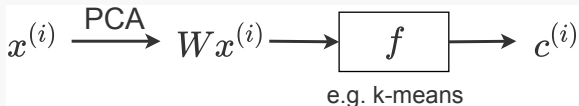
Linear PCA assumes data are separable in \mathbb{R}^n

A non-linear generalization

- ▶ Project data into higher dimension using feature mapping
 $\phi: \underline{\mathbb{R}^n} \rightarrow \underline{\mathbb{R}^d}$ ($\underline{d} \geq \underline{n}$)

Kernel PCA

Feature extraction using PCA



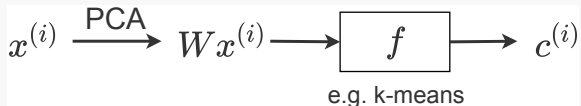
Linear PCA assumes data are separable in \mathbb{R}^n

A non-linear generalization

- ▶ Project data into higher dimension using feature mapping $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$ ($d \geq n$)
- ▶ Feature mapping is defined by a kernel function $K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^T \phi(x^{(j)})$ or kernel matrix $K \in \mathbb{R}^{m \times m}$

Kernel PCA

Feature extraction using PCA



Linear PCA assumes data are separable in \mathbb{R}^n

A non-linear generalization

- ▶ Project data into higher dimension using feature mapping $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$ ($d \geq n$)
- ▶ Feature mapping is defined by a kernel function $K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^T \phi(x^{(j)})$ or kernel matrix $K \in \mathbb{R}^{m \times m}$
- ▶ We can now perform standard PCA in the feature space

Kernel PCA

(Bernhard Schoelkopf, Alexander J. Smola, and Klaus-Robert Mueller. 1999. *Kernel principal component analysis*. In *Advances in kernel methods*) Sample covariance matrix of feature mapped data (assuming $\phi(x)$ is centered)

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \phi(x^{(i)}) \phi(x^{(i)})^T \in \mathbb{R}^{d \times d}$$

Kernel PCA

(Bernhard Schoelkopf, Alexander J. Smola, and Klaus-Robert Mueller. 1999. *Kernel principal component analysis*. In *Advances in kernel methods*) Sample covariance matrix of feature mapped data (assuming $\phi(x)$ is centered)

$$\underline{\Sigma} = \frac{1}{m} \sum_{i=1}^m \phi(x^{(i)}) \phi(x^{(i)})^T \in \mathbb{R}^{d \times d}$$

Let $(\underline{\lambda}_k, \underline{u}_k)$, $k = 1, \dots, d$ be the eigen decomposition of Σ :

$$\Sigma u_k = \lambda_k u_k$$

Kernel PCA

(Bernhard Schoelkopf, Alexander J. Smola, and Klaus-Robert Mueller. 1999. *Kernel principal component analysis*. In *Advances in kernel methods*) Sample covariance matrix of feature mapped data (assuming $\phi(x)$ is centered)

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \phi(x^{(i)}) \phi(x^{(i)})^T \in \mathbb{R}^{d \times d}$$

Let $(\lambda_k, u_k), k = 1, \dots, d$ be the eigen decomposition of Σ :

$$\Sigma u_k = \lambda_k u_k$$

PCA projection of $x^{(l)}$ onto the k th principal component u_k :

$$\underline{\underline{\phi(x^{(l)})^T u_k}}$$

Kernel PCA

(Bernhard Schoelkopf, Alexander J. Smola, and Klaus-Robert Mueller. 1999. *Kernel principal component analysis*. In *Advances in kernel methods*) Sample covariance matrix of feature mapped data (assuming $\phi(x)$ is centered)

$$\Sigma = \frac{1}{m} \sum_{i=1}^m \phi(x^{(i)}) \phi(x^{(i)})^T \in \mathbb{R}^{d \times d}$$

Let $(\lambda_k, u_k), k = 1, \dots, d$ be the eigen decomposition of Σ :

$$\Sigma u_k = \lambda_k u_k$$

PCA projection of $x^{(l)}$ onto the k th principal component u_k :

$$\underline{\phi(x^{(l)})^T u_k}$$

How to avoid evaluating $\phi(x)$ explicitly?

The Kernel Trick

Represent projection $\phi(x^{(l)})^T u_k$ using kernel function K :

- Write u_k as a linear combination of $\phi(x^{(1)}), \dots, \phi(x^{(m)})$:

$$u_k = \sum_{i=1}^m \alpha_k^i \phi(x^{(i)})$$

$$\sum u_k = \lambda_k u_k$$

$$\left(\frac{1}{m} \sum_{i=1}^m \phi(x^{(i)}) \phi(x^{(i)})^T \right) u_k = \lambda_k u_k.$$

$$\frac{1}{m} \lambda_k \sum_{i=1}^m \underbrace{\phi(x^{(i)}) \phi(x^{(i)})^T}_{\alpha_k^i} u_k = u_k.$$

$$\sum_{i=1}^m \underbrace{\frac{1}{m} \lambda_k \phi(x^{(i)})^T u_k}_{\alpha_k^i} \phi(x^{(i)}) = u_k.$$

The Kernel Trick

Represent projection $\phi(x^{(l)})^T u_k$ using kernel function K :

- ▶ Write u_k as a linear combination of $\phi(x^{(1)}), \dots, \phi(x^{(m)})$:

$$u_k = \sum_{i=1}^m \alpha_k^i \phi(x^{(i)})$$

- ▶ PCA projection of $\underline{x^{(l)}}$ using kernel function \underline{K} :

$$\underline{\phi(x^{(l)})}^T \underline{u_k} = \underline{\phi(x^{(l)})}^T \underbrace{\sum_{i=1}^m \alpha_k^i \underline{\phi(x^{(i)})}} = \sum_{i=1}^m \alpha_k^i \underline{K(x^{(l)}, x^{(i)})}$$

How to find α_k^i 's directly ?

The Kernel Trick

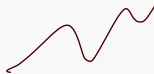
Kth eigenvector equation:

$$\sum u_k = \left(\frac{1}{m} \sum_{i=1}^m \phi(x^{(i)}) \phi(x^{(i)})^T \right) u_k = \lambda_k u_k$$

- ▶ Substitute $u_k = \sum_{i=1}^m \alpha_k^{(i)} \phi(x^{(i)})$, we obtain

$$K \alpha_k = \lambda_k m \alpha_k$$

where $\alpha_k = \begin{bmatrix} \alpha_k^1 \\ \vdots \\ \alpha_k^m \end{bmatrix}$ can be solved by eigen decomposition of K



The Kernel Trick

Kth eigenvector equation:

$$\sum u_k = \left(\frac{1}{m} \sum_{i=1}^m \phi(x^{(i)}) \phi(x^{(i)})^T \right) u_k = \lambda_k u_k$$

- ▶ Substitute $u_k = \sum_{i=1}^m \alpha_k^{(i)} \phi(x^{(i)})$, we obtain

$$K \alpha_k = \lambda_k m \alpha_k$$

where $\alpha_k = \begin{bmatrix} \alpha_k^1 \\ \vdots \\ \alpha_k^m \end{bmatrix}$ can be solved by eigen decomposition of K

- ▶ Normalize α_k such that $u_k^T u_k = 1$:

$$u_k^T u_k = \sum_{i=1}^m \sum_{j=1}^m \alpha_k^i \alpha_k^j \phi(x^{(i)})^T \phi(x^{(j)}) = \alpha_k^T K \alpha_k = \lambda_k m (\alpha_k^T \alpha_k)$$

$$\|\alpha_k\|^2 = \frac{1}{\lambda_k m}$$

Kernel PCA

When $\mathbb{E}[\phi(x)] \neq 0$, we need to center $\phi(x)$:

$$\tilde{\phi}(x^{(i)}) = \phi(x^{(i)}) - \frac{1}{m} \sum_{l=1}^m \tilde{\phi}(x^{(l)})$$

Kernel PCA

When $\mathbb{E}[\phi(x)] \neq 0$, we need to center $\phi(x)$:

$$\tilde{\phi}(x^{(i)}) = \phi(x^{(i)}) - \frac{1}{m} \sum_{l=1}^m \phi(x^{(l)})$$

The “centralized” kernel matrix is

$$\tilde{K}_{i,j} = \tilde{\phi}(x^{(i)})^T \tilde{\phi}(x^{(j)})$$

In matrix notation:

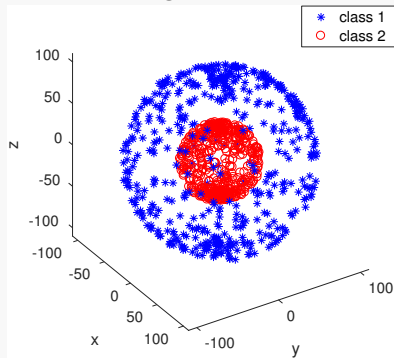
$$\tilde{K} = K - \mathbf{1}_m K - K \mathbf{1}_m + \mathbf{1}_m K \mathbf{1}_m$$

where $\mathbf{1}_m = \begin{bmatrix} 1/m & \dots & 1/m \\ \vdots & \ddots & \vdots \\ 1/m & \dots & 1/m \end{bmatrix} \in \mathbb{R}^{m \times m}$

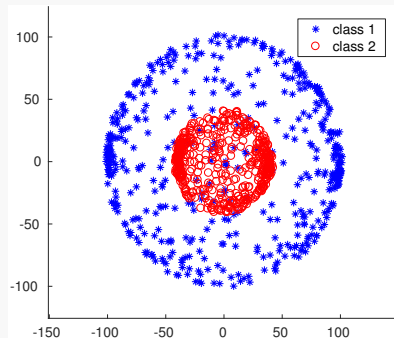
Use \tilde{K} to compute PCA

Kernel PCA Example

original data

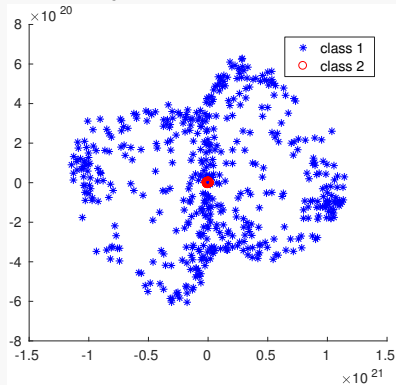


standard PCA



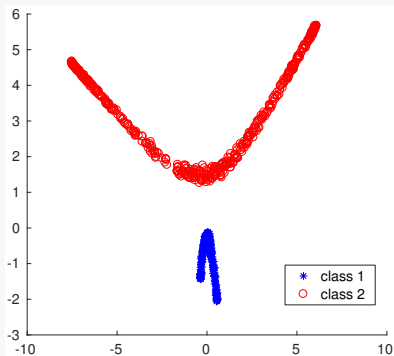
Kernel PCA Example

Polynomial kernel PCA



$$k(x, x') = (x \cdot x' + 1)^5$$

Gaussian kernel PCA



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

Discussions of kernel PCA

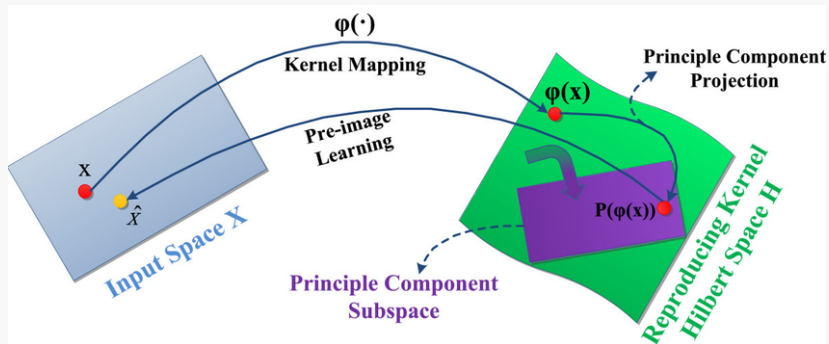
- ▶ Often used in clustering, abnormality detection, etc

Discussions of kernel PCA

- ▶ Often used in clustering, abnormality detection, etc
- ▶ Requires finding eigenvectors of $m \times m$ matrix instead of $n \times n$

Discussions of kernel PCA

- ▶ Often used in clustering, abnormality detection, etc
- ▶ Requires finding eigenvectors of $m \times m$ matrix instead of $n \times n$
- ▶ Dimension reduction by projecting to k -dimensional principal subspace is generally not possible



The Pre-Image problem: reconstruct data in input space x from feature space vectors $\phi(x)$

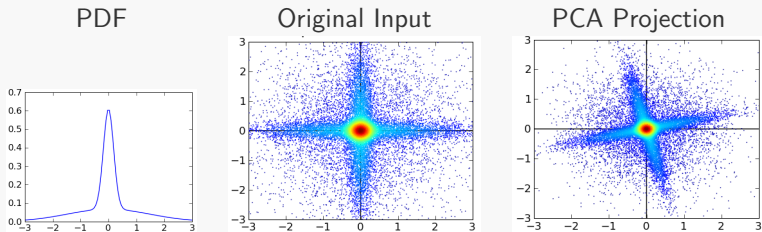
PCA Limitations

- ▶ Assumes input data is real and continuous
- ▶ Assumes **approximate normality** of input space (but may still work well on non-normally distributed data in practice) ← *sample mean & covariance must be sufficient statistics*

PCA Limitations

- ▶ Assumes input data is real and continuous
- ▶ Assumes **approximate normality** of input space (but may still work well on non-normally distributed data in practice) ← *sample mean & covariance must be sufficient statistics*

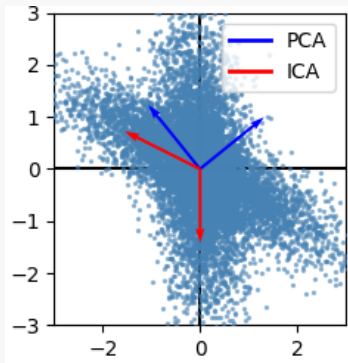
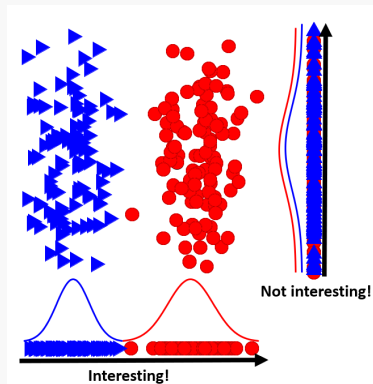
Example of strongly non-normal distributed input:



PCA Limitations

PCA results may not be useful when

- ▶ Axes of larger variance is less 'interesting' than smaller ones.
- ▶ Axes of variations are not orthogonal;



Representation learning

- ▶ Transform input features into “simpler” or “interpretable” representations.
- ▶ Used in feature extraction, dimension reduction, clustering etc

Summary

Representation learning

- ▶ Transform input features into “simpler” or “interpretable” representations.
- ▶ Used in feature extraction, dimension reduction, clustering etc

Unsupervised learning algorithms:

	low dimension	sparse	disentangle variations
k-means	✓	✓	
spectral embedding	✓		✓
PCA	✓		✓