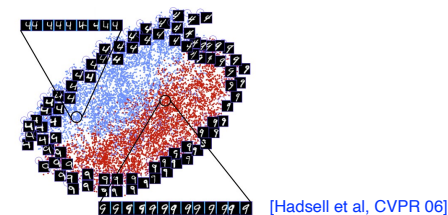


Dimensionality reduction

The slides are closely adapted from Subhransu Maji's slides

Motivation

- Data visualization
 - Hard to visualize data that lives in **high dimensions** — reduce it two **two** or **three** dimensions for visualization



- Curse of dimensionality
 - Some **learning methods** don't scale well with the **number of features** (e.g., kNN, kernel density estimators)
 - Lower **memory overhead** and **training/testing time**
 - Fewer dimensions is a form of **regularization**

2

Dimensionality reduction

- The goal is to reduce the **dimension** of the data in **high-dimensions** (say 10000) to **low dimensions** (say 2) while retaining the “**important**” characteristics of the data
- **Unsupervised setting**, so the notion of **important** characteristics is hard to define
- Closely related to **clustering**
 - **Clustering**: reduce the number of **data**
 - **Dimensionality reduction**: reduce the number of **features**



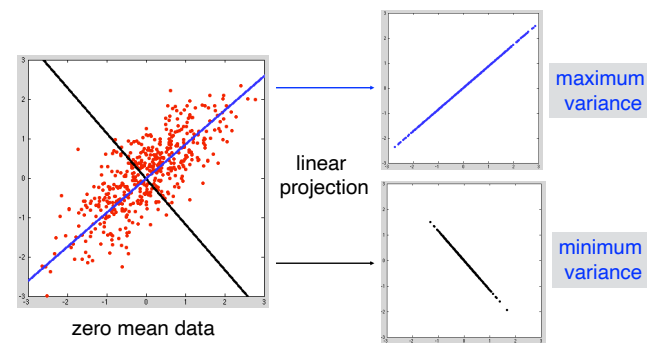
$$\mathbf{x}_i \in R^D, i = 1, 2, \dots, N$$

$$\text{data matrix} = R^{N \times D}$$

3

Linear dimensionality reduction

- All you can do is **project** the **data** onto a **vector** and use the projected distances as the **embeddings**
- **Example**: projecting two dimensional data to one



4

Optimal linear projection

- Find a linear projection that maximizes the variance of the projection
- Assume we have data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \in \mathbb{R}^D$ of zero mean
- Let \mathbf{u} be the projection vector
- Let the projections of the data p_1, p_2, \dots, p_N

- The mean of the projections is zero $p_i \leftarrow \mathbf{x}_i^T \mathbf{u}$

$$\sum_i p_i = \sum_i \mathbf{x}_i^T \mathbf{u} = \left(\sum_i \mathbf{x}_i \right)^T \mathbf{u} = 0$$

- Maximize the variance of the projection:

$$\max_{\mathbf{u}} \sum_i (\mathbf{x}_i^T \mathbf{u})^2 \text{ subject to: } \|\mathbf{u}\| = 1$$

5

Optimal linear projection

- Lets rewrite this in matrix notation
- Let \mathbf{X} be the $N \times D$ data matrix (each row is data point)
- The projection vector \mathbf{u} is a $D \times 1$ matrix
- The vector of projections is given by \mathbf{Xu} , a $N \times 1$ matrix
- We can rewrite the optimization as:

$$\max_{\mathbf{u}} \|\mathbf{Xu}\|^2 \text{ subject to: } \mathbf{u}^T \mathbf{u} = 1$$

- The corresponding Lagrangian is:

$$\mathcal{L}(\mathbf{u}, \lambda) = \|\mathbf{Xu}\|^2 - \lambda(\mathbf{u}^T \mathbf{u} - 1)$$

- At maxima:

$$\Delta_{\mathbf{u}} = 2\mathbf{X}^T \mathbf{Xu} - 2\lambda \mathbf{u}$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X}) \mathbf{u} = \lambda \mathbf{u} \quad \leftarrow \text{eigenvalue problem}$$

6

Optimal linear projection

- Compute the data covariance matrix $\mathbf{X}^T \mathbf{X}$

$$[\mathbf{X}^T \mathbf{X}]_{ij} = \sum_n \mathbf{x}_{ni} \mathbf{x}_{nj}$$

- The optimal (maximal variance) projection direction is the first eigenvector of the data covariance matrix

$$(\mathbf{X}^T \mathbf{X}) \mathbf{u} = \lambda \mathbf{u}$$

- What about learning a second projection direction?
 - For non-redundancy additionally require that $\mathbf{v}^T \mathbf{u} = 0$

$$\max_{\mathbf{v}} \|\mathbf{Xv}\|^2 \text{ subject to: } \mathbf{v}^T \mathbf{v} = 1, \mathbf{v}^T \mathbf{u} = 0$$

- This is the second eigenvector of the data covariance matrix

7

Optimal linear projection

- First find the first component \mathbf{u}
- Reconstruct the data using it \mathbf{Xu}
- Find the residual $\mathbf{X}' = \mathbf{X} - \mathbf{Xu} \left(\frac{\mathbf{u}}{\|\mathbf{u}\|} \right)^T$
- Find the next component using \mathbf{X}'
 - Since $\mathbf{X}' \mathbf{u} = 0$ the new \mathbf{u} will be orthogonal to the previous one.
 - Basically, we are removing all signals along the same component

Principal component analysis (PCA)

Find U and Λ such that

$$X^T X = U \Lambda U^T$$

where $U^T U = I$ and Λ is a diagonal matrix.

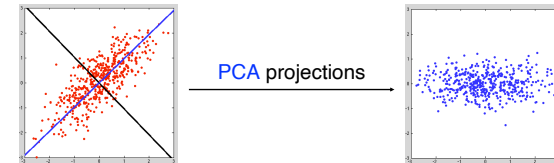
Columns of U are eigen vectors and Λ is a diagonal matrix of λ 's.

Xu is the embedding of X in the direction of u .

Principal component analysis (PCA)

- Generalizing this argument leads to [principal component analysis](#)
- The [eigenvectors](#) give you the [projection](#) directions — to compute the [embeddings](#) you have to multiply the [data](#) by the [projections](#)
- For completeness here is the [Matlab](#) code:

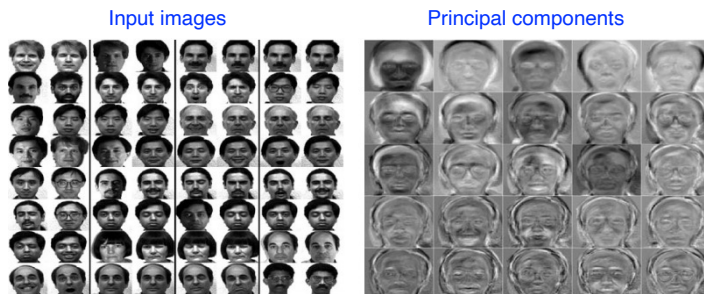
```
function [E, U, lambda] = PCA(X, K)
mu = mean(X);           % Data mean
N = size(X,1);          % Number of data
X = X - ones(N,1)*mu;   % Center the data
covX = X'*X;             % Data covariance
[U,lambda] = eigs(covX, K); % Compute top K eigenvalues
E = X*U;                 % Compute embeddings
```



10

Application: Eigenfaces

- [Eigenfaces](#) — a linear basis for face images [Turk, Pentland '91]
- Each face is a weighted linear combination of [eigenfaces](#)
- Compare faces by comparing the weights



11

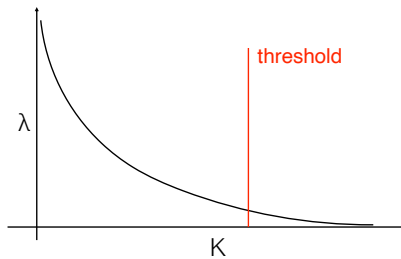
- Reconstructing face image using few components
- We need more components to see the details.



Image from: http://docs.opencv.org/2.4/modules/contrib/doc/facerec/facerec_tutorial.html

What should K be?

- For **visualization** $K = 2$ or 3
- For **dimensionality reduction** it depends on the problem
 - Option: ignore projections that correspond to small eigenvalues
 - Option: based on **computational** and **memory** constraints



13

PCA

Find U and Λ such that

$$X^T X = U \Lambda U^T$$

where $U^T U = I$ and Λ is a diagonal matrix.

Columns of U are eigen vectors and Λ is a diagonal matrix of λ 's.

- What is the dimension of covariance matrix $X^T X$ if data is n-D?
 - What if n is large, say a 100x100 image?
- It is not easy to run PCA on such a large matrix

Singular Value Decomposition(SVD)

SVD problem: Find W, V, Σ such that

$$X = W \Sigma V^T$$

Where $W^T W = I$, $V^T V = I$ and Σ is a diagonal matrix.

We can write:

$$X^T X = (W \Sigma V^T)^T (W \Sigma V^T)$$

$$X^T X = (V \Sigma W^T) (W \Sigma V^T)$$

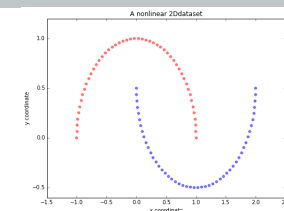
$$X^T X = V \Sigma^2 V^T$$

Compare it with PCA: $X^T X = U \Lambda U^T$

So eigen values are squares of elements in Σ

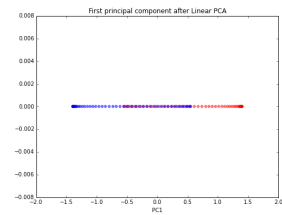
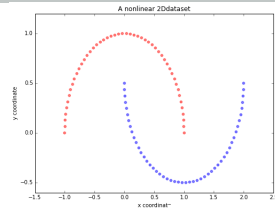
This is much faster and more stable compared to PCA for large dimensions.

Linear and kernel PCA



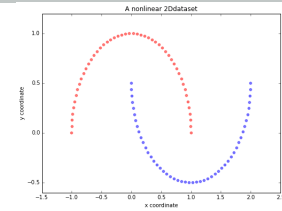
14

Linear and kernel PCA

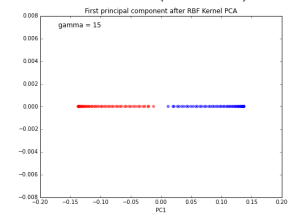
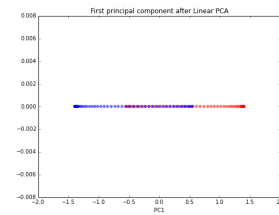


CU

Linear and kernel PCA



$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|_2^2\right)$$



21



Summary

- **Dimensionality reduction** for **visualization** or **preprocessing**
- **Linear** methods
 - **PCA** — linear projections of data
solve $(\mathbf{X}^T \mathbf{X})\mathbf{x} = \lambda \mathbf{x}$ — eigenvectors of covariance matrix
- **Non-linear** methods
 - **kernel PCA** — linear projections in kernel space
solve $\mathbf{K}\mathbf{x} = \lambda \mathbf{x}$ — eigenvectors of the kernel matrix
- **There are several methods that we didn't discuss**
 - Spectral clustering, ISOMAP, Locally linear embedding, tSNE, etc

24

Slides credit

- Most slides are adapted from Subhransu Maji's course and CIML book by Hal Daume III
- Linear and kernel PCA notes: http://pca.narod.ru/scholkopf_kernel.pdf
- The example for kernel PCA is from:
http://sebastianraschka.com/Articles/2014_kernel_pca.html

25