## Exploratory Analysis Dimensionality Reduction

#### Davide Bacciu

#### Computational Intelligence & Machine Learning Group Dipartimento di Informatica Università di Pisa bacciu@di.unipi.it

Introduzione all'Intelligenza Artificiale - A.A. 2012/2013





#### Lecture Outline



- Exploratory Analysis
- Dimensionality Reduction
  - Curse of Dimensionality
  - General View
- 3 Feature Extraction
  - Finding Linear Projections
  - Principal Component Analysis
  - Applications and Advanced Issues

Conclusion

#### Drowning into complex data



Slide credit goes to Percy Liang (Lawrence Berkeley National Laboratory)

## Exploratory Data Analysis (EDA)

#### • Discover structure in data

- Find unknown patterns in the data that cannot be predicted using current expert knowledge
- Formulate new hypotheses about the causes of the observed phenomena
- A mix of graphical and quantitative techniques
  - Visualization
  - Finding informative attributes in the data
  - Finding natural groups in the data
- Interdisciplinary approach
  - Computer graphics
  - Machine learning
  - Data Mining
  - Statistics

### A Machine Learning Perspective

#### • Often an unsupervised learning task

- Dimensionality reduction
  - Feature Extraction
  - Feature Selection
- Clustering
- Tackle with
  - Large datasets..
  - ...as well as high-dimensional data and small sample size
- Exploiting tools and models beyond statistics
  - E.g. non-parametric neural models

#### Finding Natural Groups in DNA Microarray



SL Pomeroy et al (2002) Prediction of central nervous system embryonal tumour outcome based on gene

expression, Nature, 415, 436-442

#### Finding Informative Genes



SL Pomeroy et al (2002) Prediction of central nervous system embryonal tumour outcome based on gene expression, *Nature*, 415, 436-442

Curse of Dimensionality General View

### The Curse of Dimensionality

If the data lies in a high dimensional space, then an enormous amount of data is required to learn a model

- Curse of Dimensionality (Bellman, 1961)
- Some problems become intractable as the number of the variables increases
  - Huge amount of training data required
  - Too many model parameters (complexity)



Given a fixed number of training samples, the predictive power reduces as sample dimensionality increases (Hughes Effect, 1968)

Curse of Dimensionality General View

A Simple Combinatorial Example (I)

A toy 1-dimensional classification task with 3 classes

Classes cannot be separated well: lets add another feature..





Better class separation, but still errors. What if we add another feature?

Curse of Dimensionality General View

#### A Simple Combinatorial Example (II)



$$\begin{array}{c|c} x_1 & \longrightarrow \{1,2,3\} \\ x_2 & \longrightarrow \{1,2,3\} \\ x_3 & \longrightarrow \{1,2,3\} \end{array}$$

Classes are well separated

- Exponential growth in the complexity of the learned model with increasing dimensionality
- Exponential growth in the number of examples required to maintain a given sampling density
  - 3 samples per bin in 1-D
  - 81 samples per bin in 3-D

Curse of Dimensionality General View

#### **Intrinsic Dimension**

The intrinsic dimension of data is the minimum number of independent parameters needed to account for the observed properties of the data



Data might live in a lower dimensional surface (fold) than expected

Curse of Dimensionality General View

What is the Intrinsic Dimension?

Might not be an easy question to answer...



It may increase due to noise

A data fold needs to be unfolded to reveal its intrinsic dimension

Curse of Dimensionality General View

#### Informative Vs Uninformative Features

Data can be made of several dimensions that are either unimportant or comprise only noise

- Irrelevant information might distract the learning model
- Learning resources (memory) are wasted to represent irrelevant portions of the input space

Dimensionality reduction aims at automatically finding a lower-dimensional representation of high-dimensional data

- Counteracts the curse of dimensionality
- Reduces the effect of unimportant attributes

Curse of Dimensionality General View

## Why Dimensionality Reduction?

#### Data Visualization

- Projecting high-dimensional data to a 2D/3D screen space
- Preserving topological relationships
- E.g. visualize semantically related textual documents
- Data Compression
  - Reducing storage requirements
  - Reducing complexity
  - E.g. stopwords removal
- Feature ranking and selection
  - Identifying informative bits of information
  - Noise reduction
  - E.g. identify words correlated with document topics

Curse of Dimensionality General View

#### Flavors of Dimensionality Reduction

Feature Extraction - Create a lower dimensional representation of x ∈ R<sup>D</sup> by combining the existing features with a given function f : R<sup>D</sup> → R<sup>D'</sup>

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_D \end{bmatrix} \xrightarrow{\mathbf{y} = f(\mathbf{x})} \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_{D'} \end{bmatrix}$$

• Feature Selection - Choose a *D*'-dimensional subset of all the features (possibly the most-informative)

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_D \end{bmatrix} \xrightarrow{\text{select } i_1, \dots, i_{D'}} \mathbf{y} = \begin{bmatrix} x_{i_1} \\ x_{i_2} \\ \cdots \\ x_{i_{D'}} \end{bmatrix}$$

Curse of Dimensionality General View

#### A Unique Formalization

#### **Definition (Dimensionality Reduction)**

Given an input feature space  $\mathbf{x} \in \mathbb{R}^{D}$  find a mapping  $f : \mathbb{R}^{D} \to \mathbb{R}^{D'}$  such that D' < D and  $\mathbf{y} = f(\mathbf{x})$  preserves most of the informative content in  $\mathbf{x}$ .

Often the mapping  $f(\mathbf{x})$  is chosen as a linear function  $\mathbf{y} = \mathbf{W}\mathbf{x}$ 

- y is a linear projection of x
- $\mathbf{W} \in \mathbb{R}^{D'} \times \mathbb{R}^{D}$  is the matrix of linear coefficients

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_{D'} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1D} \\ w_{21} & w_{22} & \dots & w_{2D} \\ \dots & \dots & \dots & \dots \\ w_{D'1} & w_{D'2} & \dots & w_{D'D} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_D \end{bmatrix}$$

Curse of Dimensionality General View

# Unsupervised Vs Supervised Dimensionality Reduction

The linear/nonlinear map  $\mathbf{y} = f(\mathbf{x})$  is learned from the data based on an error function that we seek to minimize

- Signal representation (Unsupervised)
  - The goal is to represent the samples accurately in a lower-dimensional space
  - Principal Component Analysis (PCA)
- Classification (Supervised)
  - The goal is to enhance the class-discriminatory information in the lower-dimensional space
  - Linear Discriminant Analysis (LDA)



Feature 1

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Feature Extraction

Objective - Create a lower dimensional representation of  $\mathbf{x} \in \mathbb{R}^{D}$  by combining the existing features with a given function  $f : \mathbb{R}^{D} \to \mathbb{R}^{D'}$ , while preserving as much information as possible

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \cdots \\ x_D \end{bmatrix} \xrightarrow{\mathbf{y} = f(\mathbf{x})} \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_{D'} \end{bmatrix}$$

where  $D' \ll D$  and, for visualization, D' = 2 or D' = 3.

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Linear Feature Extraction

- Signal Representation (Unsupervised)
  - Independent Component Analysis (ICA)
  - Principal Component Analysis (PCA)
  - Non-negative Matrix Factorization (NMF)
- Classification (Supervised)
  - Linear Discriminant Analysis (LDA)
  - Canonical Correlation Analysis (CCA)
  - Partial Least Squares (PLS)

We focus on unsupervised approaches exploiting linear mapping functions

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Linear Methods Setup

Given *N* samples  $\mathbf{x}_n \in \mathbb{R}^D$ , define the input data as the matrix

$$\mathbf{X} = \begin{bmatrix} x_{11} & | & \dots & x_{N1} \\ \dots & \mathbf{x}_2 & \dots & \dots \\ x_{1D} & | & \dots & x_{ND} \end{bmatrix} \in \mathbb{R}^D \times \mathbb{R}^N$$

Choose  $D' \ll D$  projection directions  $\mathbf{w}_k$ 

$$\mathbf{W} = \begin{bmatrix} w_{11} & | & \dots & w_{D'1} \\ \dots & \mathbf{w}_2 & \dots & \dots \\ w_{1D} & | & \dots & x_{D'D} \end{bmatrix} \in \mathbb{R}^D \times \mathbb{R}^{D'}$$

Compute the projection of **x** along each direction  $\mathbf{w}_k$  as

$$\mathbf{y} = [y_1, \ldots, y_{D'}]^T = \mathbf{W}^T \mathbf{x}$$

Linear methods only differ in the criteria used for choosing W

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

Linear Projection - A Graphical Interpretation



3D samples projected on an hyperplane generated by 2 projection directions



## 2D projection of the input samples on the hyperplane

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

## Principal Component Analysis (PCA)

Orthogonal linear projection of high dimensional data onto a low dimensional subspace preserving as much variance information as possible



Objective

- Minimize the projection error, i.e. the error of the reconstructed sample ||x<sub>n</sub> - x̃<sub>n</sub>||
- Maximize the variance of the projected data Y

The good news is that both objectives are equivalent!!

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### PCA - Two Operations

Encode Project data onto the principal components

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}$$
 for k – th component  $y_k = \mathbf{w}_k^T \mathbf{x}$ 

Decode Reconstruct the projected data

$$\tilde{\mathbf{x}} = \mathbf{W}\mathbf{y} = \sum_{k=1}^{D'} y_k \mathbf{w}_k$$

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

## PCA -Variance Maximization

Given *N* samples  $\{\mathbf{x}_n\}_{n=1}^N$  and  $\mathbf{x}_n \in \mathbb{R}^D$ 

#### Goal

Project data into a D' < D dimensional space such that the variance of the projected data is maximized

For simplicity consider D' = 1

- A single projection direction w<sub>1</sub>
- Assume normalized vectors  $\|\mathbf{w}_1\|_2 = 1$ 
  - Orthonormal basis from numerical analysis
  - Serves to select a single solution among infinite w

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

Variance Maximization - Input Space

• Compute the means of the input data  $\{\mathbf{x}_n\}_{n=1}^N$ 

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

Compute the covariance of the input data

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

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

Variance Maximization - Projected Data

- Compute the means of the projected data as  $\mathbf{w}_1^T \overline{\mathbf{x}}$
- Compute the variance of the projected data as

$$\frac{1}{N}\sum_{n=1}^{N}\left\{\mathbf{w}_{1}^{T}\mathbf{x}_{n}-\mathbf{w}_{1}^{T}\overline{\mathbf{x}}\right\}^{2}=\frac{1}{N}\sum_{n=1}^{N}\left\{\mathbf{w}_{1}^{T}(\mathbf{x}_{n}-\overline{\mathbf{x}})\right\}^{2}$$
$$=\frac{1}{N}\sum_{n=1}^{N}\mathbf{w}_{1}^{T}(\mathbf{x}_{n}-\overline{\mathbf{x}})(\mathbf{x}_{n}-\overline{\mathbf{x}})^{T}\mathbf{w}_{1}$$
$$=\mathbf{w}_{1}^{T}\mathbf{S}\mathbf{w}_{1}$$

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

PCA - Variance Maximization Problem

Goal - Maximizing the variance of the projected data

$$\mathcal{L} = \max_{\mathbf{w}} \left\{ \mathbf{w}^{\mathsf{T}} \mathbf{S} \mathbf{w} 
ight\}$$

subject to the normalization constraint

$$\|{f w}\|_2 = 1$$

How? Don't panic! No theoretical explanation. For that you will need to take the Machine Learning course

Basically, it is an optimization problem that it is solved by differentiating  $\mathcal{L}$  to find its maximum

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

PCA - Variance Maximization Solution

For D' = 1 the solution is the first principal component  $\mathbf{w} = \mathbf{u}_1$  such that

$$\mathbf{Su}_1 = \lambda_1 \mathbf{u}_1$$

where

- $\lambda_1 \in \mathbb{R}$  is the first eigenvalue of **S** (i.e. the largest)
- $\mathbf{u}_1 \in \mathbb{R}^D$  is the associate first eigenvector
- $\lambda_1$  is the variance of the projected data, i.e.

$$\lambda_1 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

Maximize the variance  $\Rightarrow$  choose eigenvector  ${\bf u}$  with largest associated eigenvalue  $\lambda$ 

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

## PCA - More Principal Components

What if I want more than 1 projection direction (D' > 1)?

- Choose each new direction **w**<sub>k</sub> as one that
  - Maximizes the variance of projected data
  - Is subject to the normalization constraint  $\|\mathbf{w}_k\|_2 = 1$
  - Is orthogonal to those already selected, i.e.  $\mathbf{w}_1, \ldots, \mathbf{w}_{k-1}$

The solution is in the eigenvectors of the input covariance S

- The covariance **S** of a *D*-dimensional input space has *D* eigenvectors
  - The eigenvector u<sub>1</sub> of the largest eigenvalue λ<sub>1</sub> is the first principal component
  - The eigenvector  $\mathbf{u}_2$  of the second-largest eigenvalue  $\lambda_2$  is the second principal component
  - The eigenvector u<sub>3</sub> of the third-largest eigenvalue λ<sub>3</sub> is the third principal component

• . . .

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

PCA Solution - Eigenvalue Decomposition

The PCA solution reduces to finding the eigenvalue decomposition of the covariance matrix of input data

#### $\mathbf{S} = \mathbf{U} \wedge \mathbf{U}^{\mathcal{T}}$

where

- $\mathbf{U} = [\mathbf{u}_k]_{k=1}^D$  is the  $D \times D$  matrix of eigenvectors  $\mathbf{u}_k$
- Λ is the *D* × *D* diagonal matrix whose diagonal element λ<sub>k</sub> is the *k*-th eigenvalue

A D' < D dimensional projection space is created by choosing D' eigenvectors  $\{\mathbf{u}_k\}_{k=1}^{D'}$  corresponding to the D' largest eigenvalues  $\{\lambda_k\}_{k=1}^{D'}$ 

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Practical PCA (I)

Step 1 Organize Data - Put your *N* samples into a  $D \times N$  matrix **X** Step 2 Compute Means - Calculate the empirical means of your data

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

Step 3 Preprocess Data - Subtract means  $\overline{\mathbf{x}}$  to each input sample

$$\overline{\mathbf{X}} = \mathbf{X} - \overline{\mathbf{x}}$$

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Practical PCA (Ia)

#### Input data Compute means Rescale data



Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Practical PCA (II)

Step 4 Compute Covariance - Calculate the covariance of input data

$$\mathbf{S} = \frac{1}{N} \overline{\mathbf{X}} \overline{\mathbf{X}}^T$$

Step 5 Eigenvalue Decomposition - Compute the eigenvalue decomposition of the covariance

 $\mathbf{S} = \mathbf{U} \wedge \mathbf{U}^T$ 

where 
$$\Lambda = diag(\lambda_1, \dots, \lambda_D)$$
 and  $\lambda_1 \ge \lambda_2 \ge \dots \ge \lambda_D$ 

Eigenvalue decomposition can be obtained using standard vector algebra or numerical routines (e.g. Singular Value Decomposition (SVD))

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Practical PCA (III)

Step 6 Model Selection - Select D' < D projection directions, associated to the first D' eigenvalues, so as to maximize the amount of variance retained in the projection

$$\mathbf{W} = \mathbf{U}_{D'} \begin{bmatrix} | & \dots & | \\ \mathbf{u}_1 & \dots & \mathbf{u}_{D'} \\ | & \dots & | \end{bmatrix} \in \mathbb{R}^D \times \mathbb{R}^{D'}$$

Step 7 Encoding - Transform the normalized data  $\overline{\mathbf{X}}$  by projecting it onto the D' principal components

$$\mathbf{Y} = \mathbf{W}^T \overline{\mathbf{X}}$$

where  $\bm{Y} \in \mathbb{R}^{D'} \times \mathbb{R}^N$  is a compressed representation of the input data

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Practical PCA (IIIa)

## Principal Components Data projected in the principal components' plane



Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### **Projecting New Data**

Given N' new input samples in  $\mathbf{X}' \in \mathbb{R}^D \times \mathbb{R}^{N'}$  they can be projected into the reduced space by

Subtracting the means

$$\overline{\mathbf{X}}' = \mathbf{X}' - \overline{\mathbf{x}}$$

Projecting onto the known principal components

$$\mathbf{Y}' = \mathbf{W}^T \overline{\mathbf{X}}'$$

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Application Example - Eigenfaces (I)

Each sample  $\mathbf{x}_n \in \mathbb{R}^D$  is a face picture with *D* pixels



The value of the *d*-th feature  $x_n(d)$  is the intensity level of the corresponding pixel



M. Turk and A. Pentland (1991) Face recognition using eigenfaces Proc. IEEE Conference on Computer Vision and

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Application Example - Eigenfaces (II)

What is a principal component? Clearly, an eigenface  $\mathbf{u}_k$ 



Eigenvectors can be shown as images depicting primitive facial features

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Application Example - Eigenfaces (III)

We can easily visualize the reconstruction of an image projected onto its eigenfaces



D' = 50 D' = 100 D' = 200

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

### How Many Principal Components?

Eigenvalues measure the fraction of variance captured by the projection



## Can be used to define a distortion measure

- Suppose we have selected *K* < *D* principal components
- The resulting distortion is

$$I = \sum_{k=K+1}^{D} \lambda_k$$

that is the proportion of variance neglected by the projection

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Is Variance so much Informative?





Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

## Linear Discriminant Analysis

Adding supervised class information into the projection function

- Linear Discriminant Analysis (LDA)
- Perform dimensionality reduction while preserving as much of the class discriminatory information as possible



- Maximum separation between the means of the projection
- Minimum variance within each projected class

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### **Nonlinear Projections**



To solve this problem either you un-fold the roll (manifold approaches) or you change the data representation (kernel methods)

Finding Linear Projections Principal Component Analysis Applications and Advanced Issues

#### Nonlinear Feature Extraction

- Signal Representation (Unsupervised)
  - Manifold learning algorithms: e.g. ISOMAP
  - Kernel Principal Component Analysis (KPCA)
- Classification (Supervised)
  - Kernel Discriminant Analysis (KDA)
  - Kernel Canonical Correlation Analysis (KCCA)

Kernels allow to use a linear model for a nonlinear problem

A kernel induces a new space by means of a non-linear mapping, where the original linear operations can be performed. *E.g. KPCA performs a linear PCA in the space created by the* 

kernel rather than in the original data space.

#### Take-home Messages

#### Exploratory data analysis

- Find new patterns in data
- Formulate new hypotheses
- Two key concepts
  - Curse of dimensionality Intractable problems
  - Intrinsic dimension Data lies in lower dimensional space
- Dimensionality Reduction
  - Feature Extraction Create new features by combining input data
  - Feature Selection Extract a subset of informative input dimension
- Linear feature extraction  $\Rightarrow$  **y** = **Wx** 
  - Models differentiate by the criteria used to chose W

## Wrapping up PCA..

- PCA is a linear transformation
  - Defined by the matrix of eigenvectors W of data covariance S
  - Preserves as much variance as possible, measured by the eigenvalues  $\boldsymbol{\Lambda}$
- A general linear transformation produces a rotation, translation and scaling of the space
  - PCA rotates the data so that is maximally decorrelated (orthonormal principal components)
- PCA is linear
  - It cannot fit well curved surfaces
  - Nonlinear models
- PCA does not account for class information
  - Supervised models (LDA)