Introduction to Statistical
Machine Learning

ⓒ2013
Christfried Webers
NICTA
The Australian National
University

ISML
2013

# Introduction to Statistical Machine Learning

## Christfried Webers

Statistical Machine Learning Group
NICTA
and
College of Engineering and Computer Science
The Australian National University

Canberra
February – June 2013

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")

*ISML*
*2013*

# Part XX

## *Principal Component Analysis*

# *Motivation*

Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
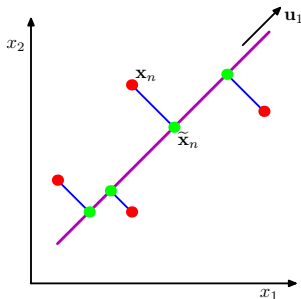University

*Motivation*

Principal Component
Analysis (PCA)

- We know already models with discrete latent variables (e.g. mixture of Gaussians), now look at continuous latent variables.
- Main goal : dimensionality reduction
- Many applications in visualisation, feature extraction, signal processing, data compression . . .
- Example: Use hand-written digits (binary data) and place them into a larger frame ($100 \times 100$) varying the position and the rotation angle.
- Data space size $= 10\,000$.
- But data live on a three-dimensional manifold ($x$, $y$, and the rotation angle).

# Principal Component Analysis (PCA)

Introduction to Statistical Machine Learning

© 2013
Christfried Webers
NICTA
The Australian National University

ISML
2013

Motivation

Principal Component Analysis (PCA)

- Idea: Linearly project the data points onto a lower dimensional subspace such that
  - the variance of the projected data is maximised, or
  - the distortion error from the projection is minimised.
- Both formulation lead to the same result.
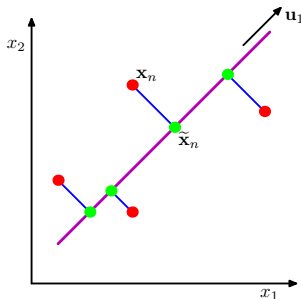- Need to find the lower dimensional subspace, called the principal subspace.

# Principal Component Analysis (PCA)

Introduction to Statistical Machine Learning

©2013
Christfried Webers
NICTA
The Australian National University

ISML
2013

- Given $N$ observations $\mathbf{x}_n \in \mathbb{R}^D$, $n = 1, \ldots, N$.
- Project onto a space with dimensionality $M < D$ while maximising the variance.
- More advanced : How to calculate $M$ from the data. Therefore here: M is fixed.
- Consider a 1-dimensional subspace spanned by some unit vector $\mathbf{u}_1 \in \mathbb{R}^D$, $\mathbf{u}_1^T \mathbf{u}_1 = 1$.

# PCA - Maximise Variance

- Each data point $\mathbf{x}_n$ is then projected onto a scalar value $\mathbf{u}_1^T \mathbf{x}_n$.
- The mean of the projected data is $\mathbf{u}_1^T \bar{\mathbf{x}}$ where $\bar{\mathbf{x}}$ is the sample mean

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

- The variance of the projected data is then

$$\frac{1}{N} \sum_{n=1}^{N} \left\{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \right\}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

with the covariance matrix

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

- Maximising $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ under the constraint $\mathbf{u}_1^T \mathbf{u}_1 = 1$ (why do we need to bound $\mathbf{u}_1$ ?) leads to the Lagrange equation

# PCA - Maximise Variance

Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
University

- Maximising $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ under the constraint $\mathbf{u}_1^T \mathbf{u}_1 = 1$ (why do we need to bound $\mathbf{u}_1$ ?) leads to the Lagrange equation

$$L(\mathbf{u}_1, \lambda_1) = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

which has a stationary point

# PCA - Maximise Variance

Motivation

Principal Component Analysis (PCA)

- Maximising $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ under the constraint $\mathbf{u}_1^T \mathbf{u}_1 = 1$ (why do we need to bound $\mathbf{u}_1$ ?) leads to the Lagrange equation

$$L(\mathbf{u}_1, \lambda_1) = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

which has a stationary point if $\mathbf{u}_1$ is an eigenvector of $\mathbf{S}$ with eigenvalue $\lambda_1$.

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1.$$

- The variance is then $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$.
- Variance is maximised if $\mathbf{u}_1$ is the eigenvector of the covariance $\mathbf{S}$ with the largest eigenvalue.

# PCA - Maximise Variance

- Continue maximising the variance amongst all possible directions orthogonal to those already considered.
- The optimal linear projection onto a $M$-dimensional space for which the variance is maximised is defined by the $M$ eigenvectors $\mathbf{u}_1, \ldots, \mathbf{u}_M$ of the covariance matrix $\mathbf{S}$ corresponding to the $M$ largest eigenvalues $\lambda_1, \ldots, \lambda_M$.
- Is this subspace always uniquely defined?

- Not if $\lambda_M = \lambda_{M+1}$.

# *PCA - Minimise Distortion Error*

- The distortion between data points $\mathbf{x}_n$ and their projection $\widetilde{\mathbf{x}}_n$

$$J = \frac{1}{N} \sum_{n=1}^{N} \|\mathbf{x}_n - \widetilde{\mathbf{x}}_n\|^2$$

  is minimised if the variance is maximised.

- The distortion error is then

$$J = \sum_{i=M+1}^{D} \lambda_i$$

  where $\lambda_i$, $i = M+1, \ldots, D$ are the smallest eigenvalues of the covariance matrix $\mathbf{S}$.
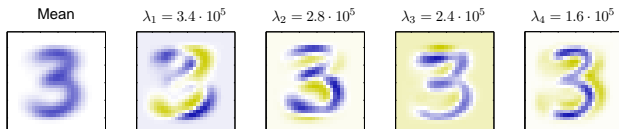
- In signal processing we speak of the signal space (principal subspace) and the noise space (orthogonal to the principal subspace).

# PCA - Applications
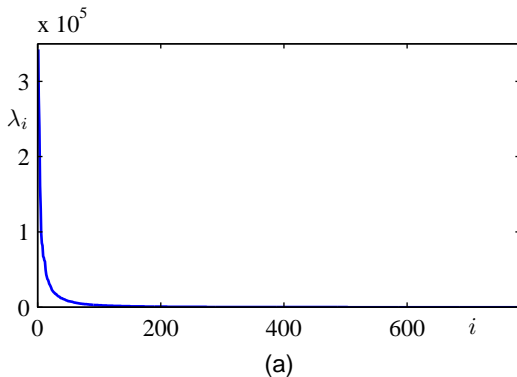
Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
University

- The eigenvectors of the covariance matrix are elements of the original vector space $u_i \in \mathbb{R}^D$.
- If the input data are images, the eigenvectors are also images.

Mean    $\lambda_1 = 3.4 \cdot 10^5$    $\lambda_2 = 2.8 \cdot 10^5$    $\lambda_3 = 2.4 \cdot 10^5$    $\lambda_4 = 1.6 \cdot 10^5$

The mean and the first four eigenvectors $\mathbf{u}_1, \ldots, \mathbf{u}_4$ of a set of handwritten digits of 'three'.
Blue corresponds to positive values, white is zero and yellow corresponds to negative values.
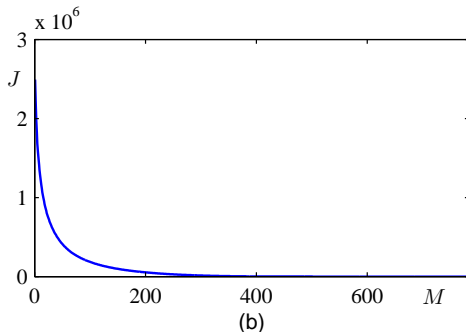
# PCA - Applications

- The eigenvalues of the covariance matrix express the variance of the data set in the direction of the corresponding eigenvectors.



Plot of the eigenvalue spectrum for the digits of three data set.

# PCA - Applications

- The sum of the eigenvalues of the covariance matrix of the discarded directions express the distortion error.

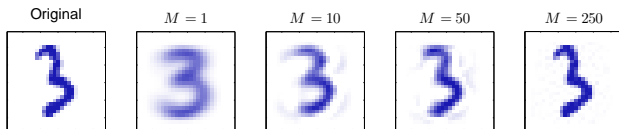$$J = \frac{1}{N} \sum_{n=1}^{N} \|\mathbf{x}_n - \widetilde{\mathbf{x}}_n\|^2$$

Plot of the distortion error versus the number of dimension of the subspace considered for projection.

# PCA - Compression

- The approximated data vector $\widetilde{\mathbf{x}}_n$ can be written in the form

$$\widetilde{\mathbf{x}}_n = \bar{\mathbf{x}} + \sum_{i=1}^{M} \left( \mathbf{u}_i^T (\mathbf{x}_n - \bar{\mathbf{x}}) \right) \mathbf{u}_i$$

- Codebook : $M + 1$ vectors of dimension $D$ ($\bar{\mathbf{x}}$ and $\mathbf{u}_i$).
- Compressed $\mathbf{x}_n$ : $M$ factors $\mathbf{u}_i^T (\mathbf{x}_n - \bar{\mathbf{x}})$

| Original | $M = 1$ | $M = 10$ | $M = 50$ | $M = 250$ |
|---|---|---|---|---|

Reconstruction of an image retaining $M$ principal components.

# PCA - Data Preprocessing

Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
University

- Standardise certain features of a data set (for instance as a preprocessing step to subsequent algorithms expecting these features).
- Usually, individual standardisation: each variable (dimension) has zero mean and unit variance. But variables are still correlated.
- PCA can do more: create decorrelated data (covariance is the identity; also called whitening or sphering of the data)
- Write the eigenvector equation for the covariance matrix $\mathbf{S}$

$$\mathbf{SU} = \mathbf{UL}$$

where $\mathbf{L}$ is the diagonal matrix of (positive!) eigenvalues.
- Transform the original data by

$$\mathbf{y}_n = \mathbf{L}^{-1/2}\,\mathbf{U}^T(\mathbf{x}_n - \bar{\mathbf{x}})$$

- The set $\{\mathbf{y}_n\}$ has mean zero and covariance given by the identity.

# PCA - Data Preprocessing

Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
University

- Transform the original data by

$$\mathbf{y}_n = \mathbf{L}^{-1/2} \, \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}})$$

- Mean of the set $\{\mathbf{y}_n\}$

$$\frac{1}{N} \sum_{n=1}^{N} \mathbf{y}_n = \frac{1}{N} \sum_{n=1}^{N} \mathbf{L}^{-1/2} \, \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}})$$

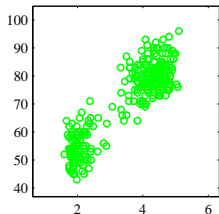$$= \mathbf{L}^{-1/2} \, \mathbf{U}^T \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_n - \bar{\mathbf{x}}) = 0$$

- Covariance of the set $\{\mathbf{y}_n\}$

$$\frac{1}{N} \sum_{n=1}^{N} \mathbf{y}_n \mathbf{y}_n^T = \frac{1}{N} \sum_{n=1}^{N} \mathbf{L}^{-1/2} \, \mathbf{U}^T (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{U} \mathbf{L}^{-1/2}$$

$$= \mathbf{L}^{-1/2} \, \mathbf{U}^T \mathbf{S} \mathbf{U} \mathbf{L}^{-1/2}$$

$$= \mathbf{L}^{-1/2} \, \mathbf{U}^T \mathbf{U} \mathbf{L} \mathbf{L}^{-1/2}$$

$$= \mathbf{I}$$

# PCA - The Effect of Whitening

Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
University

ISML
2013
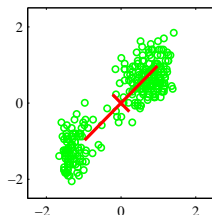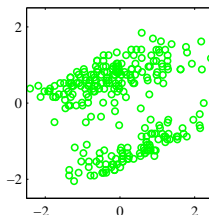
- Compare standardising and whitening of a data set.
- (b) also shows the principal axis of the normalised data set plotted as red lines over the range $\pm\lambda_i^{1/2}$.

Original data
(note the different
axis).

Standardising to
zero mean and unit
variance.

Whitening to
achieve unit
covariance.

# *Independent Component Analysis - Overview*

Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
University

- Assume we have $K$ signals and $K$ recordings, each recording containing a mixture of the signals.
- 'Cocktail party' problem : $K$ people speak at the same time in a room, and $K$ microphones pickup a mixture of what they say.
- Given unknown source signals $S \in \mathbb{R}^{N \times K}$ and an unknow mixing matrix $\mathbf{A}$, producing the observed data $X \in \mathbb{R}^{N \times K}$

$$X = SA$$

- Can we recover the original signals (Blind Source Separation)?
- Yes, under the assumption that
  - at most on of the signals is Gaussian distributed.
  - we don't care for the amplitude (including the sign).
  - we don't care for the order of the recovered signals.
  - we have at least as many observed mixtures as signals, the matrix $\mathbf{A}$ has full rank and can be inverted.

# *Independence versus Uncorrelatedness*

- Independence

$$p(x_1, x_2) = p(x_1) \, p(x_2)$$

- Uncorrelatedness (defined via a zero covariance)

$$\mathbb{E}\left[x_1 x_2\right] - \mathbb{E}\left[x_1\right] \mathbb{E}\left[x_2\right] = 0$$

- Independence implies Uncorrelatedness (prove it!).
- BUT Uncorrelatedness does NOT imply Independence.
- Example: Draw the pair $(x_1, x_2)$ with equal probability from the set $\{(0, 1), (0, -1), (1, 0), (-1, 0)\}$.
- Then $x_1$ and $x_2$ are uncorrelated because $\mathbb{E}\left[x_1\right] = \mathbb{E}\left[x_2\right] = \mathbb{E}\left[x_1 x_2\right] = 0$ .
- But $x_1$ and $x_2$ are NOT independent

$$p(x_1 = 0, x_2 = -1) = \frac{1}{4}$$
$$p(x_1 = 0) \, p(x_2 = -1) = \frac{1}{2} \times \frac{1}{4}$$

# *Independent Component Analysis - Overview*

Introduction to Statistical
Machine Learning

©2013
Christfried Webers
NICTA
The Australian National
University

- Uncorrelated variables are not necessarily independent.
- ICA maximises the statistical independence of the estimated components.
- Find $W$ (or $W^{-1}$) in such a way that

$$S = XW^{-1}$$

the columns of $S$ are maximally independent.
- Several definitions for statistical independence possible.
- Based on the concept of what is 'nongaussian'.
- Central Limit Theorem: The distribution of a sum of independent random variables tends toward a Gaussian distribution (under certain conditions).
- FastICA algorithm.