

Neural Networks and Learning Systems
TBMI 26, 2017

Lecture 6
Unsupervised Learning

Magnus Borga
magnus.borga@liu.se

Three main categories of machine learning methods

- **Supervised learning (predictive)**
Learn to generalize and classify new data based on labeled training data.
 - Pattern recognition
 - Classification
- **Unsupervised learning (descriptive)**
Discover structure and relationships in complex high-dimensional data.
- **Reinforcement learning (active)**
Generate policies/strategies that lead to a (possibly delayed) reward. Learning by doing.

Unsupervised learning

- **Task:** Find underlying structure in data.
- **Input:** Training data examples $\{x_i\}$ $i=1\dots N$.
- **Output:** Description of the data in a simpler form, e.g., with fewer dimensions or parameters.



Unsupervised learning

- Optimizes an internal cost function, e.g.
 - max variance (PCA)
 - max class separability (LDA)
- Finds a new representation of the data

Applications

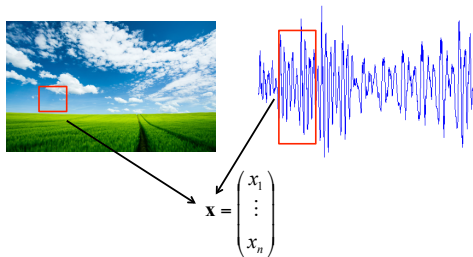
- Feature extraction
 - find order or structure in data
- Dimensionality reduction
 - keep the most “important” parts of the signal

Too many dimensions/features

Correlated features or features that do not carry any information:

- Introduce noise in the analysis/classification
- Introduce more parameters in the learning model
 - More local optima in the optimization
 - Poorer generalization
 - Higher computational effort
- Difficult to visualize high-dimensional data

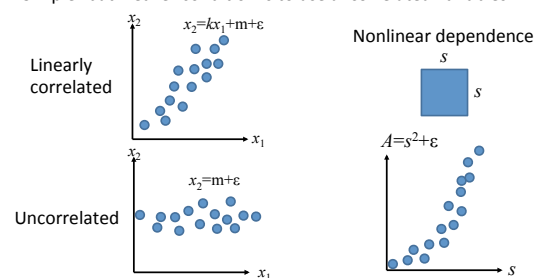
Structure in signals



Lots of structure (autocorrelation) in signals/images!
Generally **not** a good idea to use data samples directly as features.

Data/feature dependence

Want features that carry *independent* information.
Simpler but weaker condition is to use *uncorrelated* variables.



Independent vs. uncorrelated

- Statistically independent means that there is no relationship (linear or nonlinear) between variables. **Independent -> uncorrelated.**
- Uncorrelated means that there is no linear relationship between variables. **Uncorrelated does not imply independent.**
- Special case: For Gaussian distributions, uncorrelated also means independent.

Describing linear data dependence

$$\text{Var}(x) = \sigma^2 = \overset{\text{Mean value operator}}{E}[(x - \bar{x})^2] = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

$$\text{Cov}(x, y) = E[(x - \bar{x})(y - \bar{y})] = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

$$\rho(x, y) = \text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x)\text{Var}(y)}} \quad -1 \leq \rho \leq 1$$

Scaling: $\text{Var}(a \cdot x) = a^2 \text{Var}(x)$ Influenced by scaling
 $\rho(a \cdot x, b \cdot y) = \rho(x, y)$ Invariant to scaling

Multidimensional linear dependence

$$\mathbf{C} = \text{Cov}(\mathbf{x}) = E[(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T] = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Example:

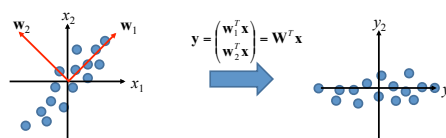
Symmetric covariance matrix

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \quad \mathbf{C} = \begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \text{Cov}(x_1, x_3) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \text{Cov}(x_2, x_3) \\ \text{Cov}(x_3, x_1) & \text{Cov}(x_3, x_2) & \text{Var}(x_3) \end{bmatrix}$$

Sometimes we also use the correlation matrix

$$\mathbf{R} = \begin{bmatrix} 1 & \text{Corr}(x_1, x_2) & \text{Corr}(x_1, x_3) \\ \text{Corr}(x_2, x_1) & 1 & \text{Corr}(x_2, x_3) \\ \text{Corr}(x_3, x_1) & \text{Corr}(x_3, x_2) & 1 \end{bmatrix}$$

Projection/transformation



Can "uncorrelate" data through a linear transformation!!

$$\mathbf{C}_x = \begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) \end{bmatrix}$$

$$\mathbf{C}_y = \begin{bmatrix} \text{Var}(y_1) & 0 \\ 0 & \text{Var}(y_2) \end{bmatrix}$$

PCA

Principal Component Analysis

- Pearson 1901
(A.k.a. Hotelling-transform or Karhunen-Loève-transform)
- Coordinate transformation to an orthogonal basis where the data is uncorrelated.
- Dimensionality reduction that preserves maximum variance (minimizes the mean square error).

Dimensionality reduction with least mean square error

- Chose a new basis whose dimensionality is smaller than the original
- The basis vectors that preserve maximum variance gives the least mean square error.

Dimensionality reduction with least mean square error

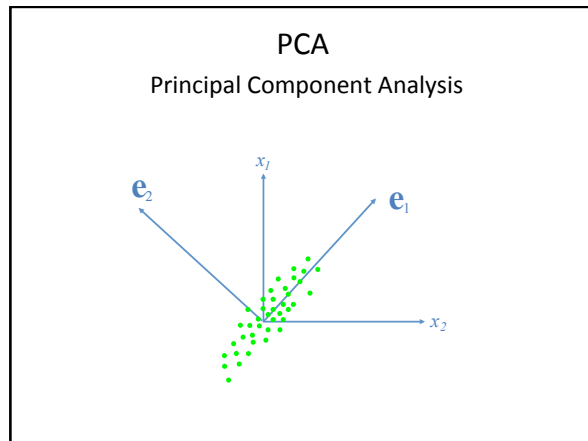
Suppose \mathbf{y} is a representation of \mathbf{x} where some dimensions are removed, i.e. some $y_i = 0$. (Assume mean = 0)

$$\varepsilon = E \left[\sum_i (x_i - y_i)^2 \right] = E \left[\sum_k x_k^2 \right] = \sum_k E[x_k^2] = \sum_k \sigma_k^2$$

k is the indices of the removed components, i.e. $k : y_k = 0$

Dimensionality reduction with least mean square error

- Expected mean square error = sum of the variance in the removed dimensions.
- The mean square error is minimized by keeping the components with highest variance.



Maximize the variance

The variance in direction $\hat{\mathbf{w}}$: (Suppose \mathbf{x} has mean 0.)

$$\sigma_{\hat{\mathbf{w}}}^2 = E[(\mathbf{x}^T \hat{\mathbf{w}})^2] = E[(\hat{\mathbf{w}}^T \mathbf{x})(\mathbf{x}^T \hat{\mathbf{w}})]$$

$$= \hat{\mathbf{w}}^T E[\mathbf{x}\mathbf{x}^T] \hat{\mathbf{w}} = \hat{\mathbf{w}}^T \mathbf{C} \hat{\mathbf{w}} = \frac{\mathbf{w}^T \mathbf{C} \mathbf{w}}{\mathbf{w}^T \mathbf{w}}$$

The covariance matrix of \mathbf{x} .

Maximize the variance

$$\sigma_{\hat{\mathbf{w}}}^2 = \frac{\mathbf{w}^T \mathbf{C} \mathbf{w}}{\mathbf{w}^T \mathbf{w}}$$

$$\frac{\partial \sigma_{\hat{\mathbf{w}}}^2}{\partial \mathbf{w}} = \frac{2}{\mathbf{w}^T \mathbf{w}} (\mathbf{C} \mathbf{w} - \sigma_{\hat{\mathbf{w}}}^2 \mathbf{w}) = 0 \Rightarrow$$

$$\mathbf{C} \mathbf{w} = \sigma_{\hat{\mathbf{w}}}^2 \mathbf{w}$$

PCA is the Eigen-value decomposition of the data covariance matrix.

The Eigen-value decomposition

$$\mathbf{C} \mathbf{e} = \lambda \mathbf{e}$$

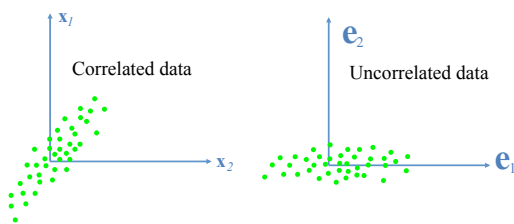
$$\mathbf{C} = \sum_i \lambda_i \mathbf{e}_i \mathbf{e}_i^T$$

If \mathbf{C} is a covariance matrix:

$$\mathbf{C} = E[(\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T] \quad \lambda_i = \sigma_{\mathbf{e}_i}^2$$

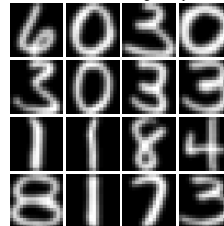
Uncorrelated components

$$E[\mathbf{e}_i^T \mathbf{x} \mathbf{x}^T \mathbf{e}_j] = \mathbf{e}_i^T \mathbf{C} \mathbf{e}_j = \lambda_j \mathbf{e}_i^T \mathbf{e}_j = 0$$



PCA - Example

About 9000 training examples

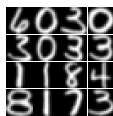


Feature vectors

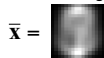
$$\begin{matrix} 16 & \text{[digit 8]} & \Rightarrow & \mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_{256} \end{pmatrix} \\ & 16 & & \end{matrix}$$

US Postal Service Digit Data
<http://www.gaussianprocess.org/gpml/data/>

PCA – Example, cont.

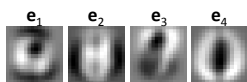


Mean digit

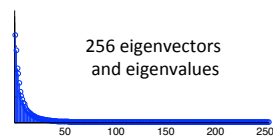
 $\bar{\mathbf{x}}$

$$\mathbf{C} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

256x256 matrix

Eigendecomposition of \mathbf{C} :

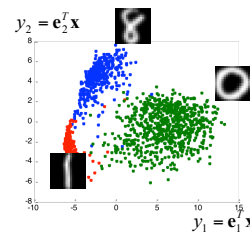
256 eigenvectors and eigenvalues



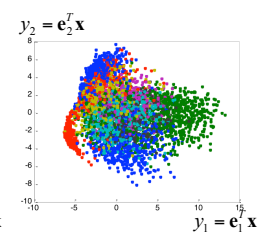
PCA – Example, cont.

From 256 to 2 feature dimensions!

3 classes

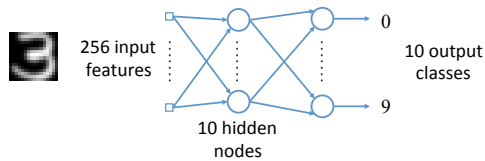


All 10 classes



PCA – Example, cont.

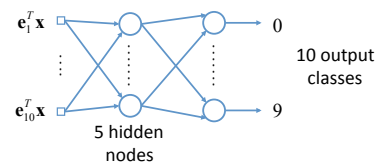
Say we try to classify with a neural network with 10 hidden nodes



parameters in network: $256 \times 10 + 10 \times 10 = 2700$
(ignoring the bias weights)

PCA – Example, cont.

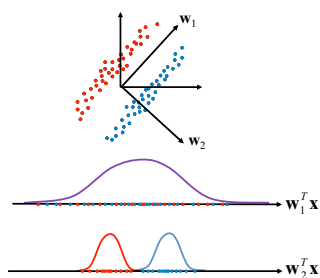
If we reduce the input dimensionality first, we may be able to do the classification with a smaller network, e.g., 10 principal components as input and 5 hidden nodes.



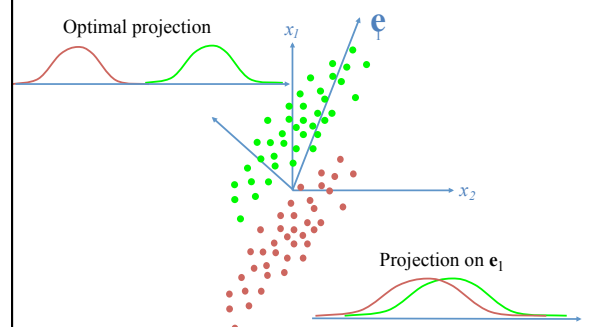
parameters in network: $10 \times 5 + 5 \times 10 = 100$
(ignoring the bias weights)

Limitations with PCA

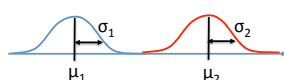
Variance is not always the most important goal!



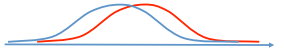
Optimal projection for separation of two clusters



Class separability



- Small variance
- Large distance



- Large variance
- Small distance

Goal: minimize variance and maximize distance.

Linear Discriminant Analysis (LDA)

a.k.a. Fishers Linear Discriminant (FLD)

- Minimize variance
- Maximize distance

$$\text{Maximize: } \varepsilon(\mathbf{w}) = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$$

LDA – Cost function $\varepsilon = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$

Distance:

$$\mu(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathbf{w}^T \mathbf{x}_i = \mathbf{w}^T \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right) = \mathbf{w}^T \bar{\mathbf{x}}$$



$$(\mu_1(\mathbf{w}) - \mu_2(\mathbf{w}))^2 = (\mathbf{w}^T (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2))^2 = \mathbf{w}^T (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{w} = \mathbf{w}^T \mathbf{M} \mathbf{w}$$

Variance:

$$\sigma^2(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T (\mathbf{x}_i - \bar{\mathbf{x}}))^2 = \dots = \mathbf{w}^T \left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \right) \mathbf{w} = \mathbf{w}^T \mathbf{C} \mathbf{w}$$

$$\sigma_1^2(\mathbf{w}) + \sigma_2^2(\mathbf{w}) = \mathbf{w}^T \mathbf{C}_1 \mathbf{w} + \mathbf{w}^T \mathbf{C}_2 \mathbf{w} = \mathbf{w}^T \mathbf{C}_{tot} \mathbf{w}$$

Exercise:
Complete all the steps!

LDA – Solution

$$\varepsilon(\mathbf{w}) = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2} = \frac{\mathbf{w}^T \mathbf{M} \mathbf{w}}{\mathbf{w}^T \mathbf{C}_{tot} \mathbf{w}}$$

Compare with PCA!

This form is called a Rayleigh quotient, which is maximized by the largest eigenvector to the generalized eigenvalue problem $\mathbf{C}_{tot} \mathbf{w} = \lambda \mathbf{M} \mathbf{w}$!

$$\text{Simplification: } \mathbf{M} \mathbf{w} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{w} = K (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

Some scalar K

$$\mathbf{w} \sim \mathbf{C}_{tot}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

Scaling of \mathbf{w} not important!

