



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

Probabilistic Data Mining

Lehel Csató

Faculty of Mathematics and Informatics
Babeş–Bolyai University, Cluj-Napoca,

November 2010



Outline

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- 1 **Modelling Data**
 - Motivation
 - Machine Learning
 - Latent variable models
- 2 Estimation methods
- 3 Unsupervised Methods



Motivation for Data Mining

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

Data Mining is **not**:

- SQL and relational data-base application;
- Storage technologies;
- Cloud Computing;

Data mining:

- The extraction of **knowledge** or **information** from an **ever-growing** collection of data.
- “Advanced” search capability that enables one to extract **patterns** useful in providing models for:
 - 1 characterising;
 - 2 prediction, and
 - 3 exploiting the data.



Data mining applications

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Identifying targets for vouchers/frequent flier bonuses or in telecommunications.
- “Basket analysis” – correlation-based analysis leading to recommending new items – Amazon.com.
- (semi)automated fraud/virus detection: use guards that protect against procedural or other types of misuse of a system.
- Forecasting e.g. energy consumption of a region for optimising coal/hydro-plants or planning;
- Exploiting textual databases – the Google business:
 - to answer user queries;
 - to put content-sensitive ads: Google AdSense



Data mining applications

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Identifying targets for vouchers/frequent flier bonuses or in telecommunications.
- “Basket analysis” – correlation–based analysis leading to recommending new items – Amazon.com.
- (semi)automated fraud/virus detection: use guards that protect against procedural or other types of misuse of a system.
- Forecasting *e.g.* energy consumption of a region for optimising coal/hydro-plants or planning;
- Exploiting textual databases – the Google business:
 - to answer user queries;
 - to put content-sensitive ads: Google AdSense



Data mining applications

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Identifying targets for vouchers/frequent flier bonuses or in telecommunications.
- “Basket analysis” – correlation–based analysis leading to recommending new items – Amazon.com.
- (semi)automated fraud/virus detection: use guards that protect against procedural or other types of misuse of a system.
- Forecasting *e.g.* energy consumption of a region for optimising coal/hydro-plants or planning;
- Exploiting textual databases – the Google business:
 - to answer user queries;
 - to put content-sensitive ads: Google AdSense



Data mining applications

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Identifying targets for vouchers/frequent flier bonuses or in telecommunications.
- “Basket analysis” – correlation–based analysis leading to recommending new items – Amazon.com.
- (semi)automated fraud/virus detection: use guards that protect against procedural or other types of misuse of a system.
- Forecasting *e.g.* energy consumption of a region for optimising coal/hydro-plants or planning;
- Exploiting textual databases – the Google business:
 - to answer user queries;
 - to put content-sensitive ads: Google AdSense



Data mining applications

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Identifying targets for vouchers/frequent flier bonuses or in telecommunications.
- “Basket analysis” – correlation–based analysis leading to recommending new items – Amazon.com.
- (semi)automated fraud/virus detection: use guards that protect against procedural or other types of misuse of a system.
- Forecasting *e.g.* energy consumption of a region for optimising coal/hydro-plants or planning;
- Exploiting textual databases – the Google business:
 - to answer user **queries**;
 - to put **content-sensitive** ads: Google AdSense



The need for data mining

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- *“Computers have promised us a fountain of wisdom but delivered a flood of data.”*
“The amount of information in the world doubles every 20 months.”
(Frawley, Piatetsky-Shapiro, Matheus, 1991)
- A competitive market environment requires sophisticated – and useful – algorithms.
- Data acquisition and storage is ubiquitous. Algorithms are required to exploit them.
- The algorithms that exploit the data-rich environment are coming usually from **the machine learning** domain.



Machine learning

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

Historical background / Motivation:

- Huge amount of **data**, that should **automatically** be processed,
- Mathematics provides general solutions, solutions are i.e. **not for a given problem**,
- Need for “science”, that uses mathematics machinery for solving **practical** problems.



Definitions for Machine Learning

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

Machine learning

Collection of methods (from statistics, probability theory) to solve problems **met in practice**.

- noise filtering for
 - non-linear regression and/or
 - non-Gaussian noise
- Classification:
 - binary,
 - multiclass,
 - partially labelled
- Clustering,
- Inversion problems,
- density estimation, novelty detection.

Generally, we need to **model the data**,



Modelling Data

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

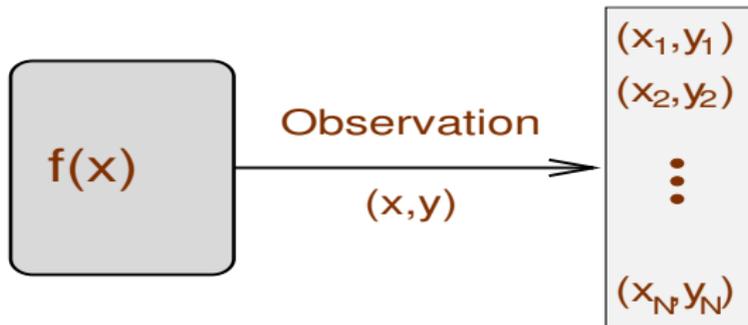
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



- Real world: there “is” a function $y = f(x)$
- Observation process: a **corrupted** datum is collected for a sample x_n :

$$t_n = y_n + \epsilon \quad \text{additive noise}$$

$$t_n = h(y_n, \epsilon) \quad h \text{ distortion function}$$

- **Problem:** find function $y = f(x)$



Latent variable models

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

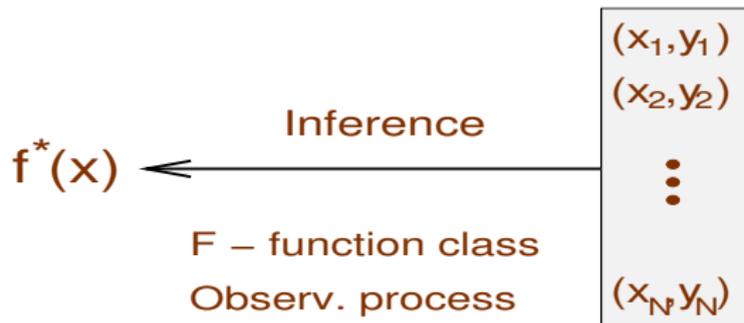
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



- **Data set** – collected.
- Assume a function class.
 - polynomial,
 - Fourier expansion,
 - Wavelet;
- Observation process – **encodes** the noise;
- Find the optimal function from the class.



Latent variable models II

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- We have the **data set** $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$.

- Consider a function class:

$$(1) \quad \mathcal{F} = \{ \mathbf{w}^T \mathbf{x} + b \mid \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \}$$

$$(2) \quad \mathcal{F} = \left\{ a_0 + \sum_{k=1}^K a_k \sin(2\pi kx) + \sum_{k=1}^K b_k \cos(2\pi kx) \right.$$

$$\left. \mid \mathbf{a}, \mathbf{b} \in \mathbb{R}^K, a_0 \in \mathbb{R} \right\}$$

- Assume an observation process:

$$y_n = f(\mathbf{x}_n) + \epsilon \quad \text{with } \epsilon \sim N(0, \sigma^2).$$



Latent variable models III

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

1 The **data set**: $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$.

2 Assume a function class:

$$\mathcal{F} = \{f(\mathbf{x}, \boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \mathbb{R}^P\}$$

\mathcal{F} – polynomial, etc.

3 Assume an observation process. Define a **loss function**:

$$L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta}))$$

For the Gaussian noise:

$$L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta})) = (y_n - f(\mathbf{x}_n, \boldsymbol{\theta}))^2.$$



Outline

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

1 Modelling Data

2 Estimation methods

- Maximum Likelihood
- Maximum a-posteriori
- Bayesian Estimation

3 Unsupervised Methods



Parameter estimation

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

Estimating parameters:

Finding the **optimal value to θ** :

$$\theta^* = \arg \min_{\theta \in \Omega} L(\mathcal{D}, \theta)$$

where

- Ω is the domain of the parameters.
- $L(\mathcal{D}, \theta)$ is a “loss function” for the data set.

Example:

$$L(\mathcal{D}, \theta) = \sum_{n=1}^N L(y_n, f(\mathbf{x}_n, \theta))$$



Maximum Likelihood Estimation

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

$L(\mathcal{D}, \boldsymbol{\theta})$ – (log)likelihood function.

Maximum likelihood estimation of the model:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \Omega} L(\mathcal{D}, \boldsymbol{\theta})$$

Example – quadratic regression:

$$L(\mathcal{D}, \boldsymbol{\theta}) = \sum_{n=1}^N (y_n - f(\mathbf{x}_n, \boldsymbol{\theta}))^2 \quad \text{– factorisation}$$

Drawback: can produce **perfect** fit to the data – **over-fitting**.



Example of an ML estimate

Graphic

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

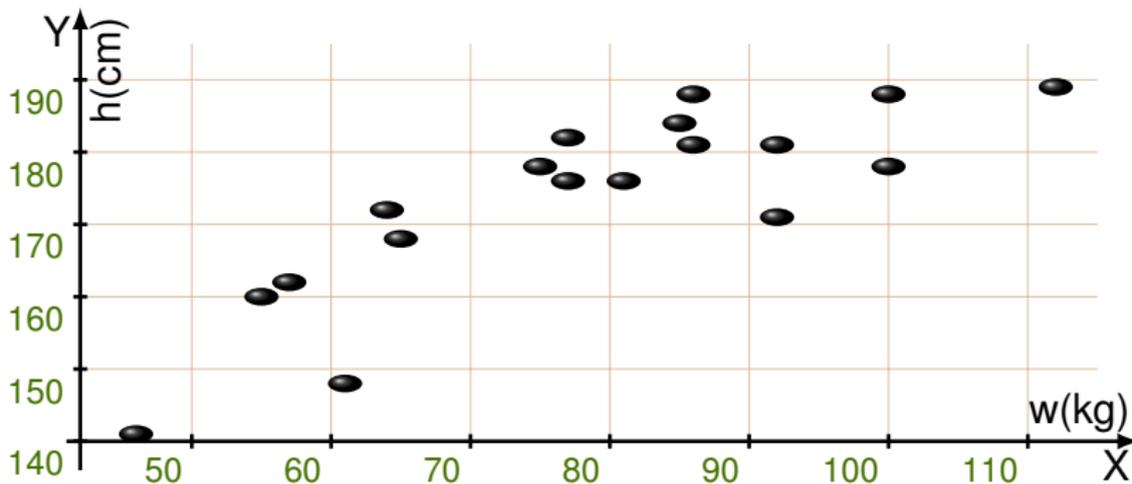
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



● We want to fit a **model** to the data.

● Use **linear model**: $h = \theta_0 + \theta_1 w$.

● Use **log-linear model**: $h = \theta_0 + \theta_1 \log(w)$.

● Use **higher order polynomials**, e.g. :

$$h = \theta_0 + \theta_1 w + \theta_2 w^2 + \theta_3 w^3 + \dots$$



Example of an ML estimate

Graphic

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

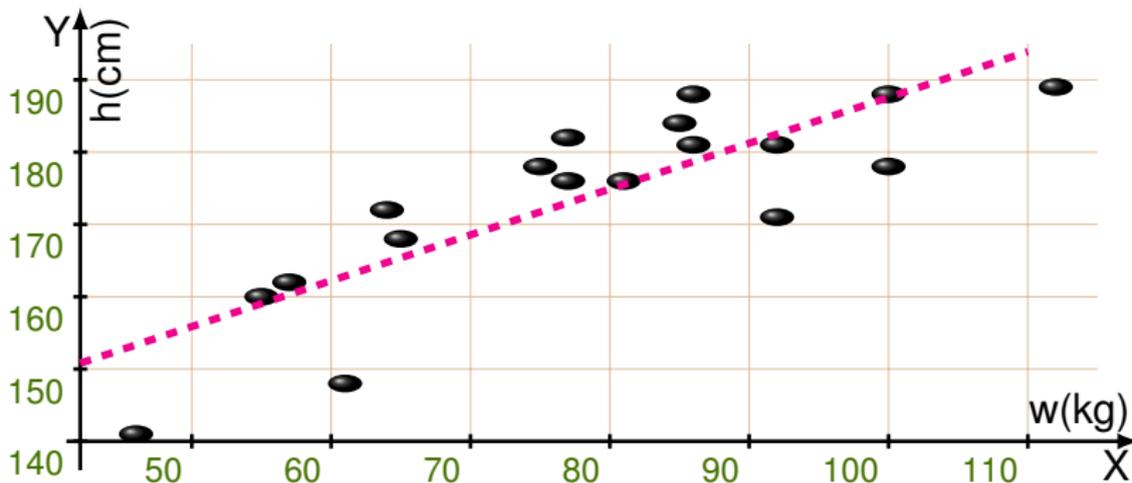
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



- We want to fit a **model** to the data.
- Use **linear model**: $h = \theta_0 + \theta_1 w$.
- Use **log-linear model**: $h = \theta_0 + \theta_1 \log(w)$.
- Use **higher order polynomials**, e.g. :

$$h = \theta_0 + \theta_1 w + \theta_2 w^2 + \theta_3 w^3 + \dots$$



Example of an ML estimate

Graphic

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

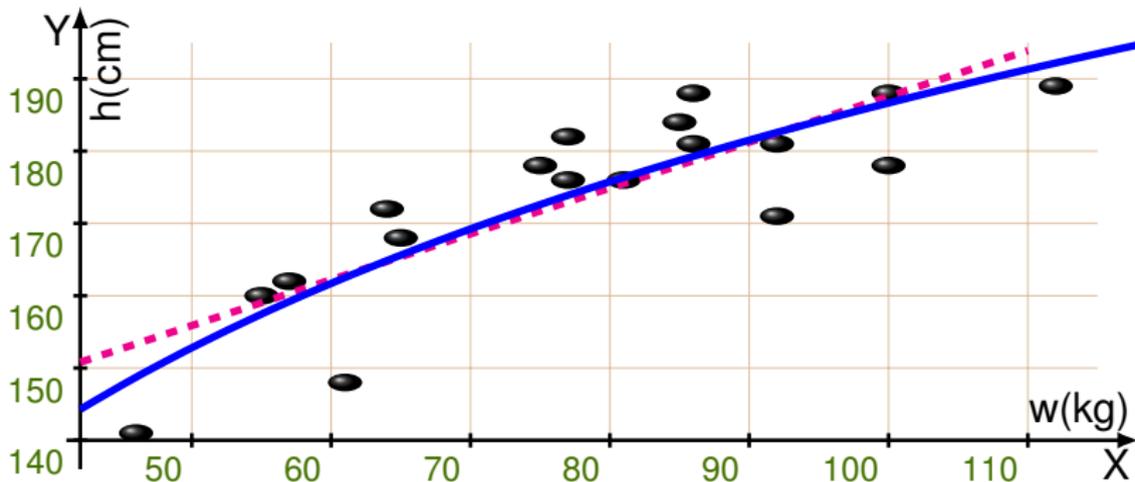
Unsupervised

General concepts

Principal Components

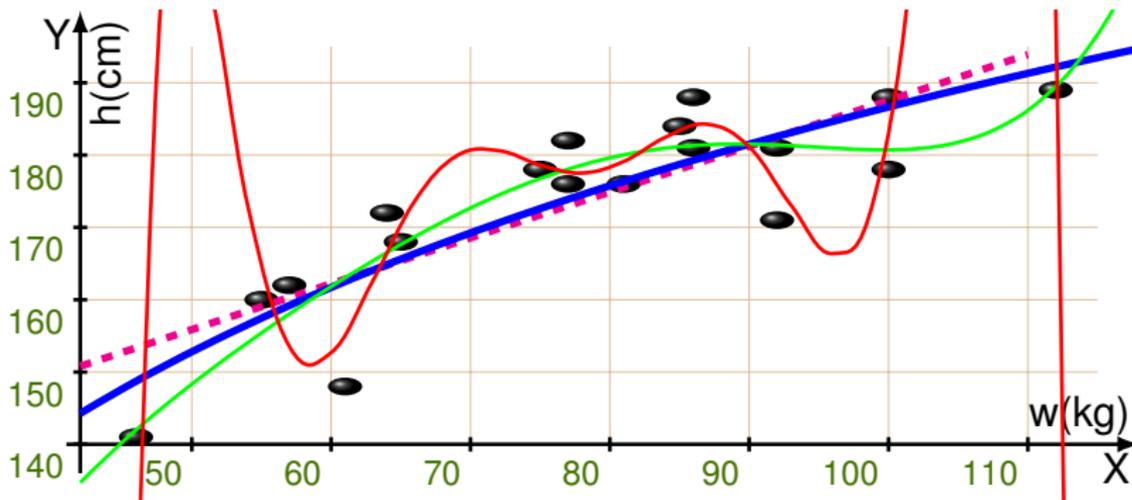
Independent Components

Mixture Models



- We want to fit a **model** to the data.
- Use **linear model**: $h = \theta_0 + \theta_1 w$.
- Use **log-linear model**: $h = \theta_0 + \theta_1 \log(w)$.
- Use **higher order polynomials**, e.g. :

$$h = \theta_0 + \theta_1 w + \theta_2 w^2 + \theta_3 w^3 + \dots$$



- We want to fit a **model** to the data.
- Use **linear model**: $h = \theta_0 + \theta_1 w$.
- Use **log-linear model**: $h = \theta_0 + \theta_1 \log(w)$.
- Use **higher order polynomials**, e.g. :

$$h = \theta_0 + \theta_1 w + \theta_2 w^2 + \theta_3 w^3 + \dots$$



Assume:

linear model for the $\mathbf{x} \rightarrow y$ relation

$$f(\mathbf{x}_n|\boldsymbol{\theta}) = \sum_{\ell=1}^d \theta_{\ell} x_{\ell}$$

with $\mathbf{x} = [1, x, x^2, \log(x), \dots]^T$

quadratic loss for $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, h_N)\}$

$$E_2(\mathcal{D}|f) = \sum_{n=1}^N (y_n - f(\mathbf{x}_n|\boldsymbol{\theta}))^2$$



Minimisation:

$$\begin{aligned}\sum_{n=1}^N (y_n - f(\mathbf{x}_n|\boldsymbol{\theta}))^2 &= (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}\end{aligned}$$

Solution:

$$\begin{aligned}0 &= 2\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - 2\mathbf{X}^T \mathbf{y} \\ \boldsymbol{\theta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

where $\mathbf{y} = [y_1, \dots, y_N]^T$ and $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T$ are the transformed data.



Generalised linear models:

- Use a *set of functions* $\Phi = [\phi_1(\cdot), \dots, \phi_M(\cdot)]$.
- Project the inputs into the space spanned by $\text{Im}(\Phi)$.
- Have a parameter vector of length M :
 $\theta = [\theta_1, \dots, \theta_M]^T$.
- The model is $\left\{ \sum_m \theta_m \phi_m(\mathbf{x}) \mid \theta_m \in \mathbb{R} \right\}$.
- The optimal parameter vector is:

$$\theta^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$



- There are many candidate **model families**:
 - the degree of polynomials specifies a model family;
 - the rank of a Fourier expansion;
 - the mixture of $\{\log, \sin, \cos, \dots\}$ also a *family*;
- Selecting the “best family” is a difficult *modelling* problem.
- In maximum likelihood there is no control on how good a family is when processing a *given* data-set.

Smaller number of parameters than $\sqrt{\#data}$.



Maximum a-posteriori

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Generalised linear model powerful – it can be extremely complex;
 - With no complexity control, overfitting problem.
- **Aim:** to include **knowledge** in the inference process.
- Our beliefs are reflected by the choice of the candidate functions.



Maximum a-posteriori

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Generalised linear model powerful – it can be extremely complex;
 - With no complexity control, overfitting problem.
- **Aim:** to include **knowledge** in the inference process.
- Our beliefs are reflected by the choice of the candidate functions.

Goal:

- Prior knowledge specification using probabilities;
- Using probability theory for consistent estimation;
- Encode the observation noise in the model;



Probabilistic data description:

- **How likely** is that θ generated the data:

$$\begin{aligned}y &= f(\mathbf{x}) & \Leftrightarrow & y - f(\mathbf{x}) \sim \delta_0 \\y &= f(\mathbf{x}) + \epsilon & \Leftrightarrow & y - f(\mathbf{x}) \sim N_\epsilon\end{aligned}$$

- **Gaussian noise:** $y - f(\mathbf{x}) \sim N(0, \sigma^2)$

$$P(y|f(\mathbf{x})) = \frac{1}{\sqrt{2\pi\sigma}} \exp \left[-\frac{(y - f(\mathbf{x}))^2}{2\sigma^2} \right]$$



William of **Ockham** (1285–1349) principle

Entities should not be multiplied beyond necessity.

Also known as (wiki...):

“Principle of simplicity” – KISS,

“When you hear hoofbeats, think horses, not zebras”.

Simple models \approx small number of parameters.

L_0 norm

L_2 norm

\Leftarrow

Probabilistic representation:

$$p_0(\boldsymbol{\theta}) \propto \exp \left[-\frac{\|\boldsymbol{\theta}\|_2^2}{2\sigma_0^2} \right]$$



M.A.P. – **probabilities** assigned to

- \mathcal{D} – via the log-likelihood function:

$$P(y_n | \mathbf{x}_n, \boldsymbol{\theta}, \mathcal{F}) \propto \exp[-L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta}))]$$

- $\boldsymbol{\theta}$ – prior probabilities:

$$p_0(\boldsymbol{\theta}) \propto \exp\left[-\frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2}\right]$$

- **A**-posteriori probability:

$$p(\boldsymbol{\theta} | \mathcal{D}, \mathcal{F}) = \frac{P(\mathcal{D} | \boldsymbol{\theta}) p_0(\boldsymbol{\theta})}{p(\mathcal{D} | \mathcal{F})}$$

$p(\mathcal{D} | \mathcal{F})$ – probability of the **data** for a given family.



M.A.P. – **probabilities** assigned to

- \mathcal{D} – via the log-likelihood function:

$$P(y_n | \mathbf{x}_n, \boldsymbol{\theta}, \mathcal{F}) \propto \exp[-L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta}))]$$

- $\boldsymbol{\theta}$ – prior probabilities:

$$p_0(\boldsymbol{\theta}) \propto \exp\left[-\frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2}\right]$$

- **A**-posteriori probability:

$$p(\boldsymbol{\theta} | \mathcal{D}, \mathcal{F}) = \frac{P(\mathcal{D} | \boldsymbol{\theta}) p_0(\boldsymbol{\theta})}{p(\mathcal{D} | \mathcal{F})}$$

$p(\mathcal{D} | \mathcal{F})$ – probability of the **data** for a given family.



M.A.P. – **probabilities** assigned to

- \mathcal{D} – via the log-likelihood function:

$$P(y_n | \mathbf{x}_n, \boldsymbol{\theta}, \mathcal{F}) \propto \exp[-L(y_n, f(\mathbf{x}_n, \boldsymbol{\theta}))]$$

- $\boldsymbol{\theta}$ – prior probabilities:

$$p_0(\boldsymbol{\theta}) \propto \exp\left[-\frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2}\right]$$

- **A-posteriori probability:**

$$p(\boldsymbol{\theta} | \mathcal{D}, \mathcal{F}) = \frac{P(\mathcal{D} | \boldsymbol{\theta}) p_0(\boldsymbol{\theta})}{p(\mathcal{D} | \mathcal{F})}$$

$p(\mathcal{D} | \mathcal{F})$ – probability of the **data** for a given family.



M.A.P. estimation – finds θ with largest probability:

$$\theta_{MAP}^* = \arg \max_{\theta \in \Omega} p(\theta | \mathcal{D}, \mathcal{F})$$

Example: with $L(y_n, f(\mathbf{x}_n, \theta))$ and Gaussian prior:

$$\theta_{MAP}^* = \operatorname{argmax}_{\theta \in \Omega} K - \frac{1}{2} \sum_n L(y_n, f(\mathbf{x}_n, \theta)) - \frac{\|\theta\|^2}{2\sigma_0^2}$$

$$\sigma_0^2 = \infty \quad \implies \quad \text{maximum likelihood.}$$

after a change of sign and $\max \rightarrow \min$



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

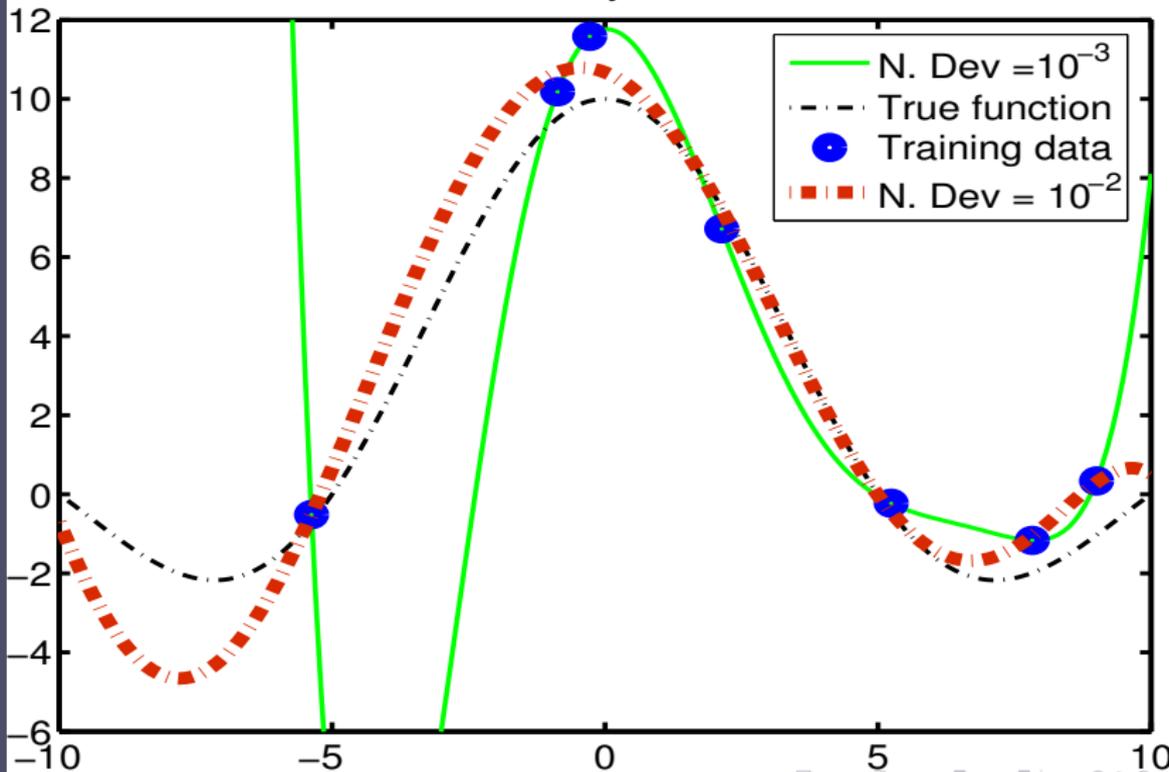
General concepts

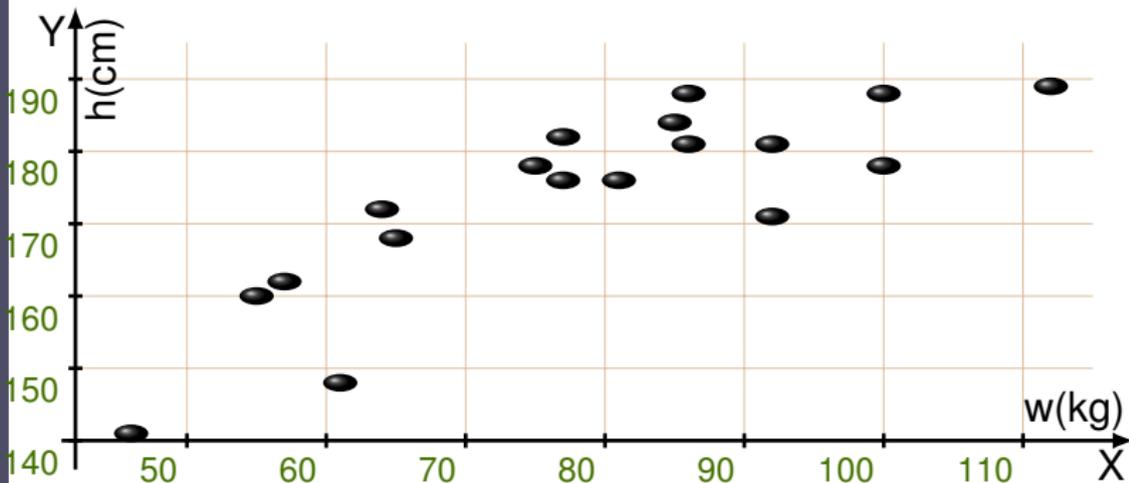
Principal Components

Independent Components

Mixture Models

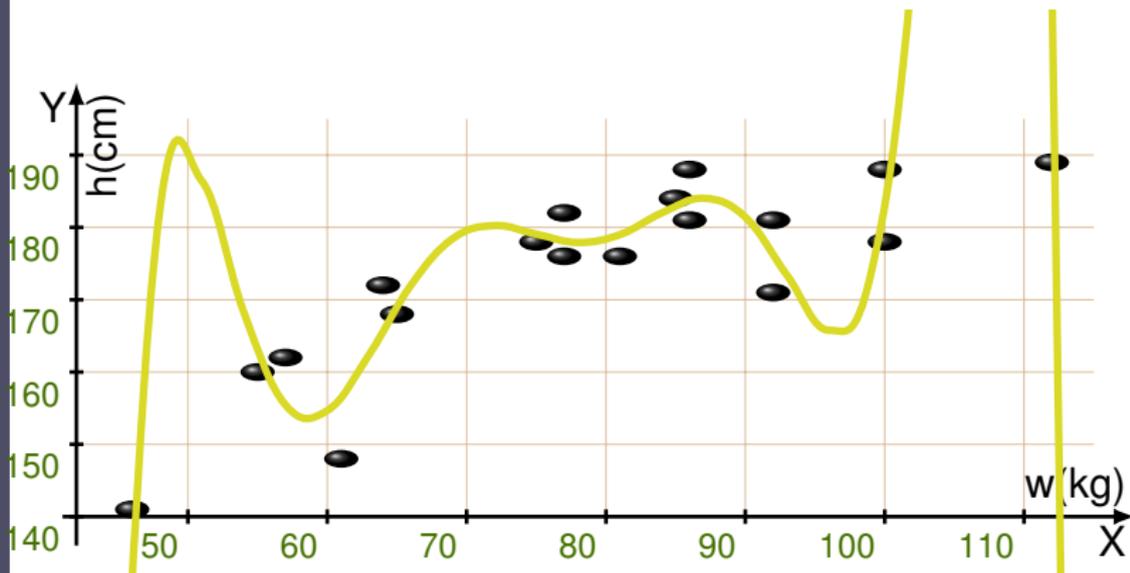
Poly 6





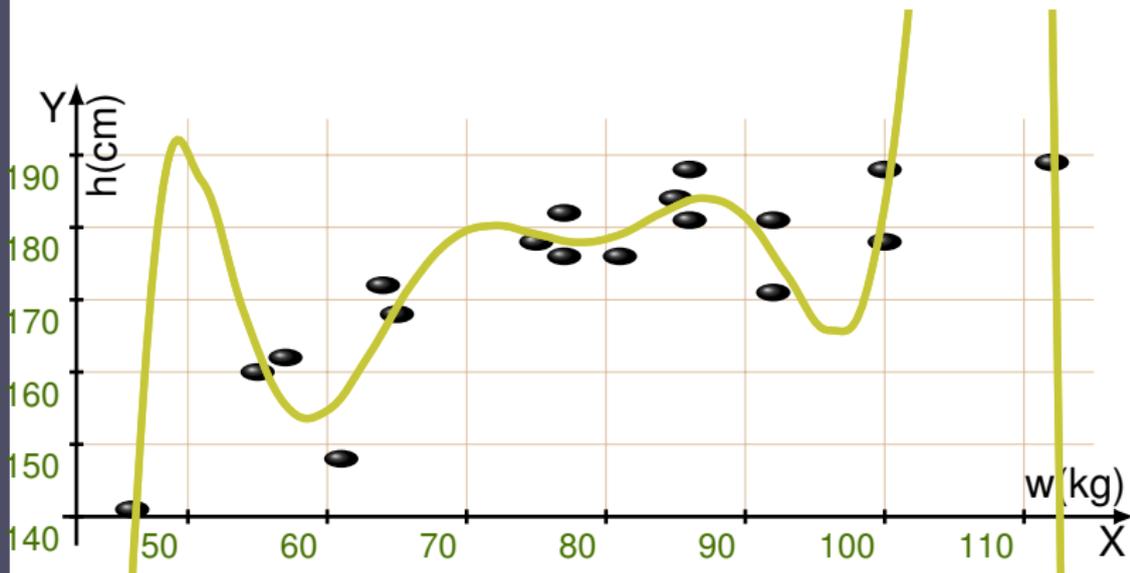
Aim: test different levels of flexibility. $\Rightarrow p = 10$

Prior width:

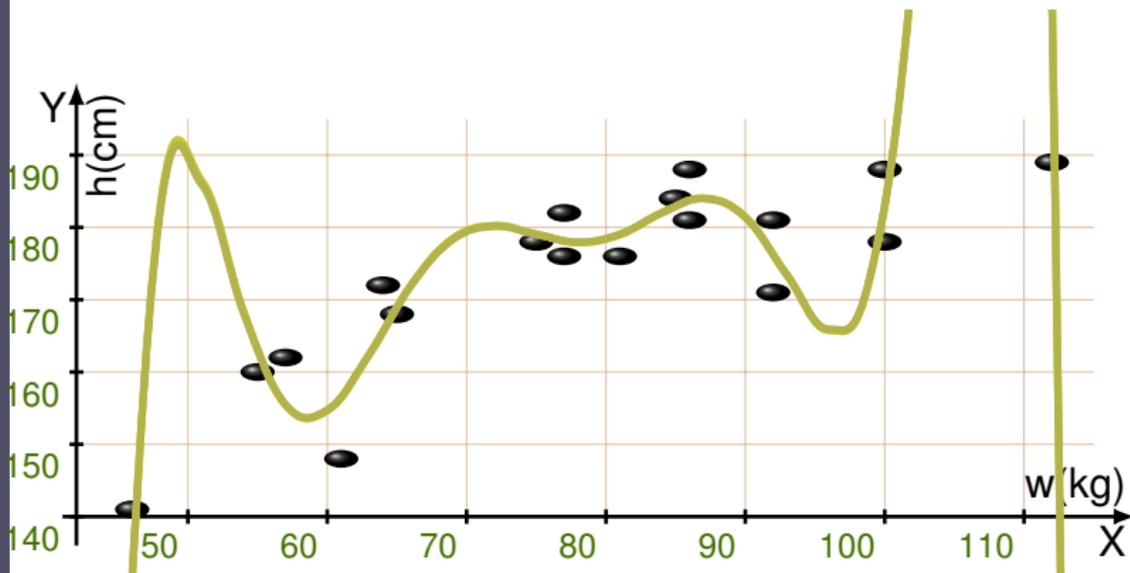


Aim: test different levels of flexibility. $\Rightarrow p = 10$

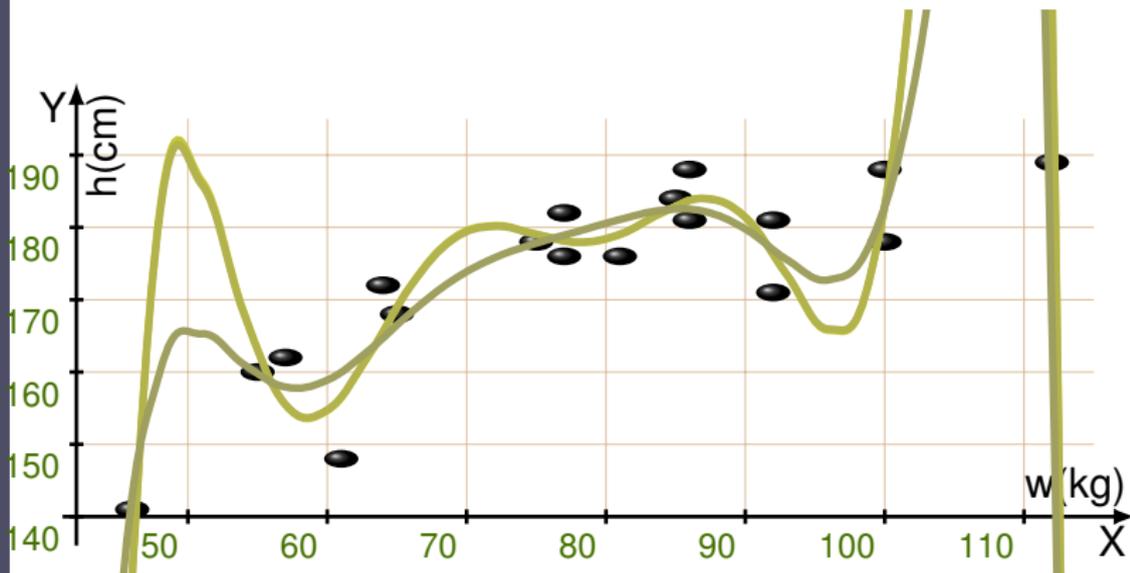
Prior width: $\sigma_0^2 = 10^6$



Aim: test different levels of flexibility. $\Rightarrow p = 10$
Prior width: $\sigma_0^2 = 10^6$ $\sigma_0^2 = 10^5$

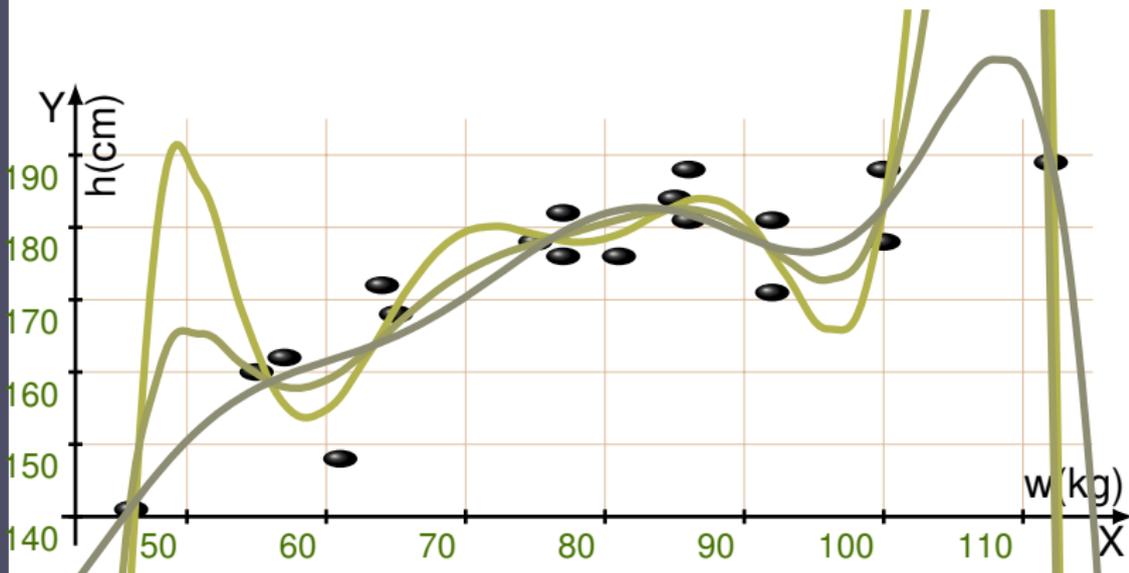


Aim: test different levels of flexibility. $\Rightarrow p = 10$
Prior width: $\sigma_0^2 = 10^6$ $\sigma_0^2 = 10^5$ $\sigma_0^2 = 10^4$



Aim: test different levels of flexibility. $\Rightarrow p = 10$

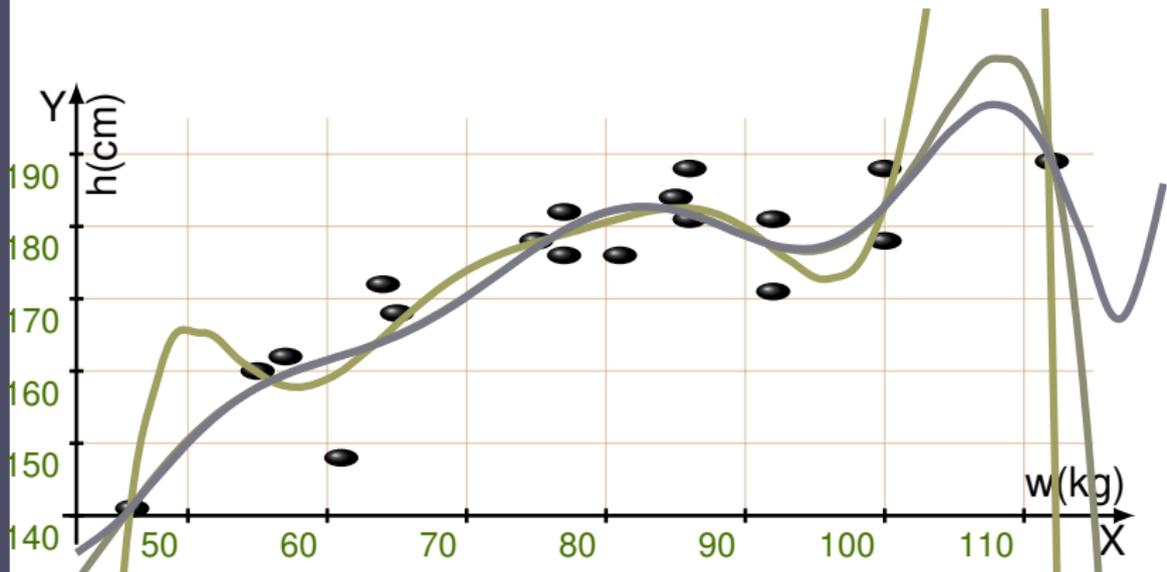
Prior width: $\sigma_0^2 = 10^6$ $\sigma_0^2 = 10^5$ $\sigma_0^2 = 10^4$ $\sigma_0^2 = 10^3$



Aim: test different levels of flexibility. $\Rightarrow p = 10$

Prior width: $\sigma_0^2 = 10^6$ $\sigma_0^2 = 10^5$ $\sigma_0^2 = 10^4$ $\sigma_0^2 = 10^3$

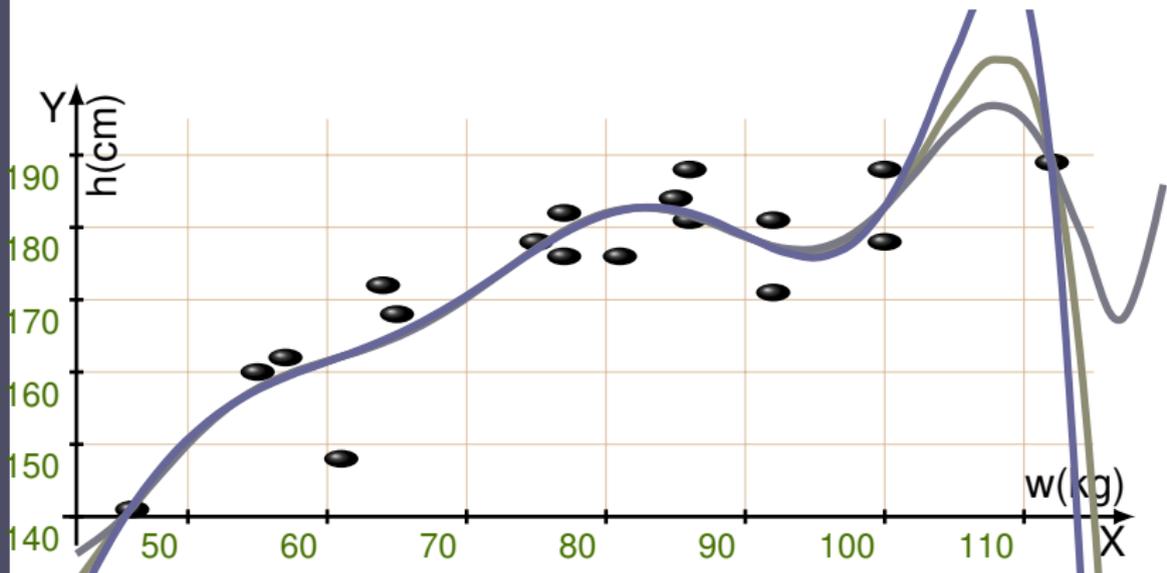
$\sigma_0^2 = 10^2$



Aim: test different levels of flexibility. $\Rightarrow p = 10$

Prior width: $\sigma_0^2 = 10^6$ $\sigma_0^2 = 10^5$ $\sigma_0^2 = 10^4$ $\sigma_0^2 = 10^3$

$\sigma_0^2 = 10^2$ $\sigma_0^2 = 10^1$



Aim: test different levels of flexibility. $\Rightarrow p = 10$

Prior width: $\sigma_0^2 = 10^6$ $\sigma_0^2 = 10^5$ $\sigma_0^2 = 10^4$ $\sigma_0^2 = 10^3$

$\sigma_0^2 = 10^2$ $\sigma_0^2 = 10^1$ $\sigma_0^2 = 10^0$



$$\boldsymbol{\theta}_{MAP}^* = \operatorname{argmax}_{\boldsymbol{\theta} \in \Omega} K - \frac{1}{2} \sum_n E_2(y_n, f(\mathbf{x}_n, \boldsymbol{\theta})) - \frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2}$$

Transform into vector notation:

$$\boldsymbol{\theta}_{MAP}^* = \operatorname{argmax}_{\boldsymbol{\theta} \in \Omega} K - \frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) - \frac{\boldsymbol{\theta}^T \boldsymbol{\theta}}{2\sigma_0^2}$$

solve for $\boldsymbol{\theta}$ by differentiation:

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) - \frac{1}{\sigma_0^2} I_d \boldsymbol{\theta} = 0$$

$$\boldsymbol{\theta}_{MAP}^* = \left(\mathbf{X}^T \mathbf{X} + \frac{1}{\sigma_0^2} I_d \right)^{-1} \mathbf{X}^T \mathbf{y}$$



Maximum a-posteriori models:

- Allow for the inclusion of prior knowledge;
- May protect against overfitting;
- Can measure the fitness of the family to the data;
Procedure called **M.L. type II**.



Idea: instead of computing the most probable value of θ , we can measure the **fit** of the model \mathcal{F} to the data \mathcal{D} .

$$\begin{aligned} P(\mathcal{D}|\mathcal{F}) &= \sum_{\theta \in \Omega} p(\mathcal{D}, \theta | \mathcal{F}) \\ &= \sum_{\theta \in \Omega} p(\mathcal{D} | \theta, \mathcal{F}) p_0(\theta | \mathcal{F}) \end{aligned}$$

Gaussian noise case and polynomial of order K :

$$\begin{aligned} \log(P(\mathcal{D}|\mathcal{F})) &= \log \left(\int_{\Omega_{\theta}} d\theta \frac{p(\mathcal{D} | \theta, \mathcal{F}) p_0(\theta | \mathcal{F})}{p(\mathcal{D} | \mathcal{F})} \right) = \log(N(\mathbf{y} | 0, \Sigma_{\mathbf{X}})) \\ &= -\frac{1}{2} \left(N \log(2\pi) + \log |\Sigma_{\mathbf{X}}| + \mathbf{y}^T \Sigma_{\mathbf{X}}^{-1} \mathbf{y} \right) \end{aligned}$$

where

$$\Sigma_{\mathbf{X}} = I_N \sigma_n^2 + \mathbf{X} \Sigma_0 \mathbf{X}^T \quad \text{with} \quad \mathbf{X} = [\mathbf{x}^0, \mathbf{x}^1, \dots, \mathbf{x}^K]$$

$$\Sigma_0 = \text{diag}(\sigma_0^2, \sigma_1^2, \dots, \sigma_K^2) = \sigma_p^2 I_{K+1}$$



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

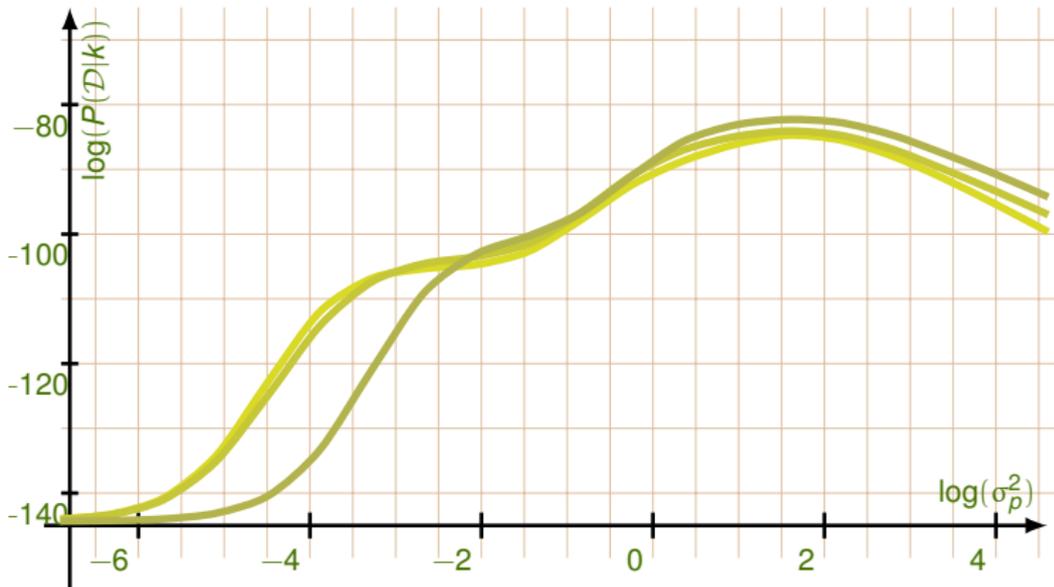
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



Aim: test different models.

Polynomial families: $k = 10$ $k = 9$ $k = 8$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

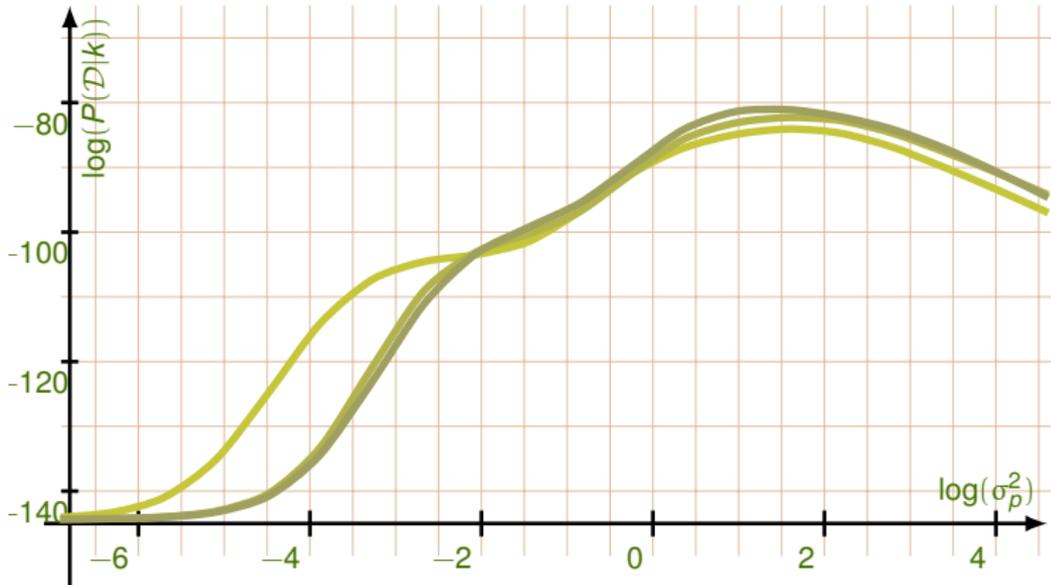
- Motivation
- Machine Learning
- Latent variable models

Estimation

- Maximum Likelihood
- Maximum a-posteriori
- Bayesian Estimation

Unsupervised

- General concepts
- Principal Components
- Independent Components
- Mixture Models



Aim: test different models.

Polynomial families: $k = 9$ $k = 8$ $k = 7$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

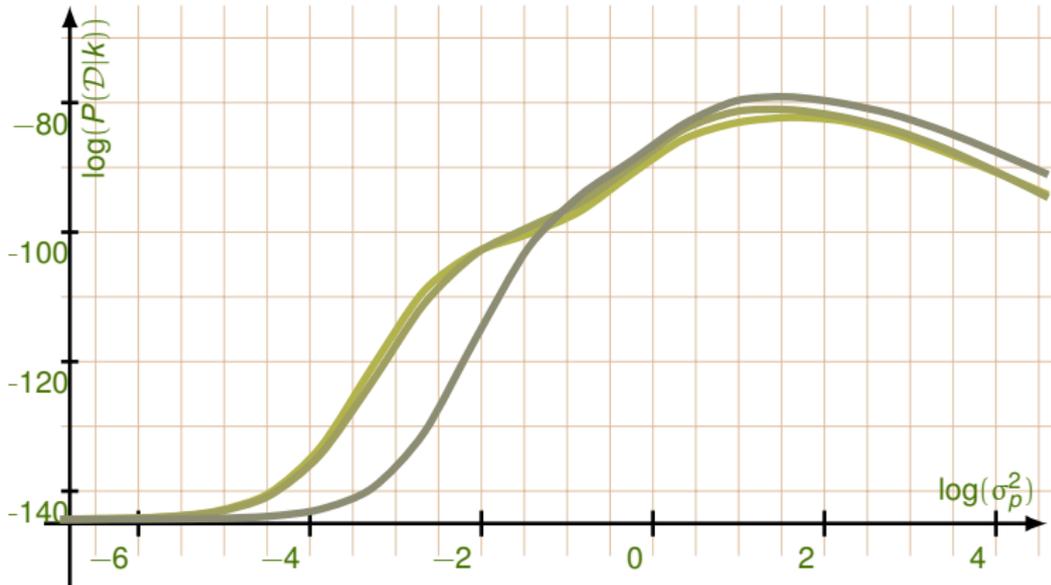
- Motivation
- Machine Learning
- Latent variable models

Estimation

- Maximum Likelihood
- Maximum a-posteriori
- Bayesian Estimation

Unsupervised

- General concepts
- Principal Components
- Independent Components
- Mixture Models



Aim: test different models.

Polynomial families: $k = 8$ $k = 7$ $k = 6$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

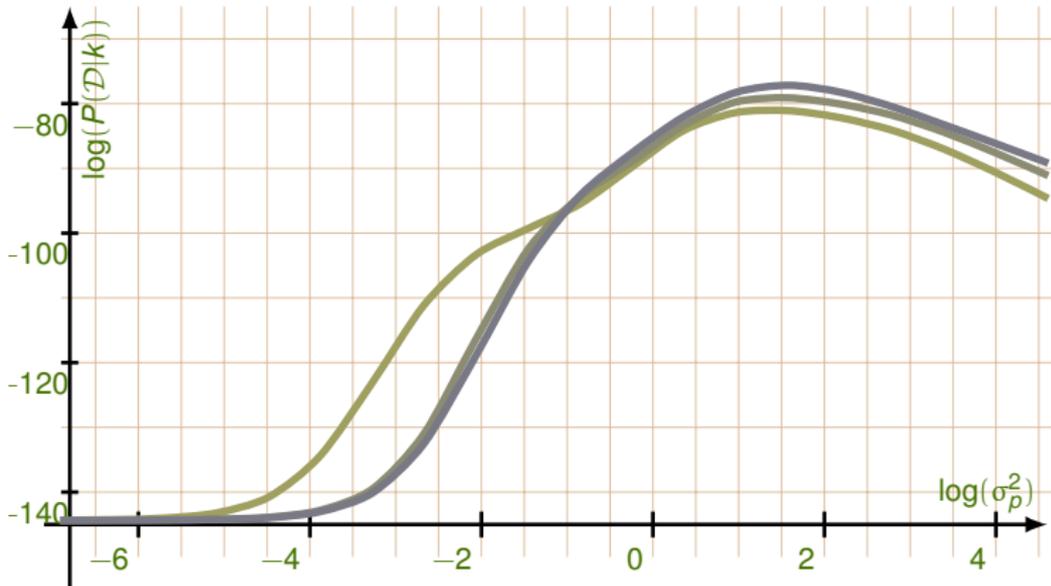
- Motivation
- Machine Learning
- Latent variable models

Estimation

- Maximum Likelihood
- Maximum a-posteriori
- Bayesian Estimation

Unsupervised

- General concepts
- Principal Components
- Independent Components
- Mixture Models



Aim: test different models.

Polynomial families: $k = 7$ $k = 6$ $k = 5$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

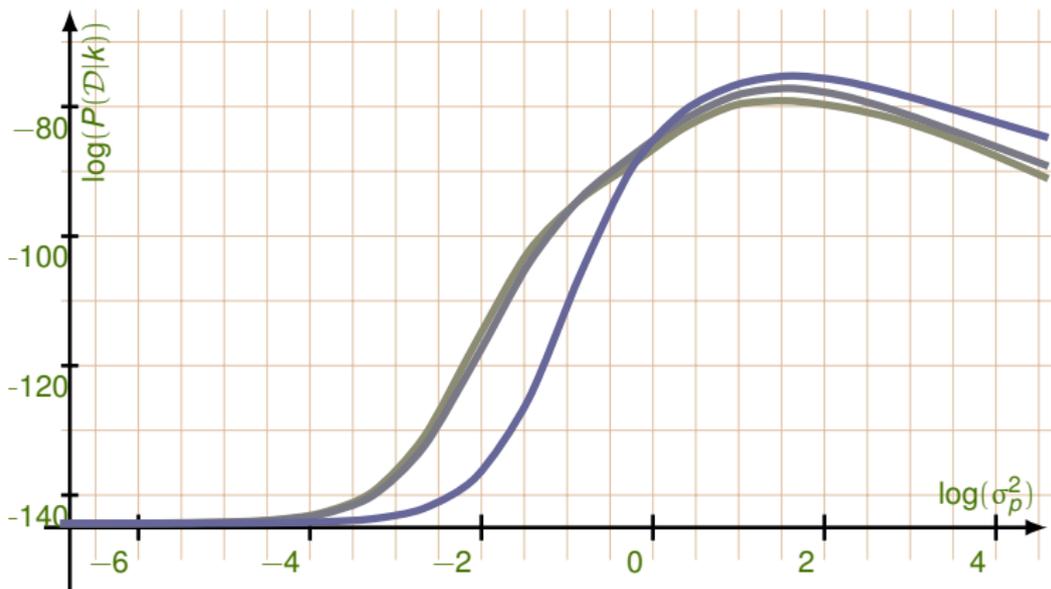
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



Aim: test different models.

Polynomial families: $k = 6$ $k = 5$ $k = 4$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

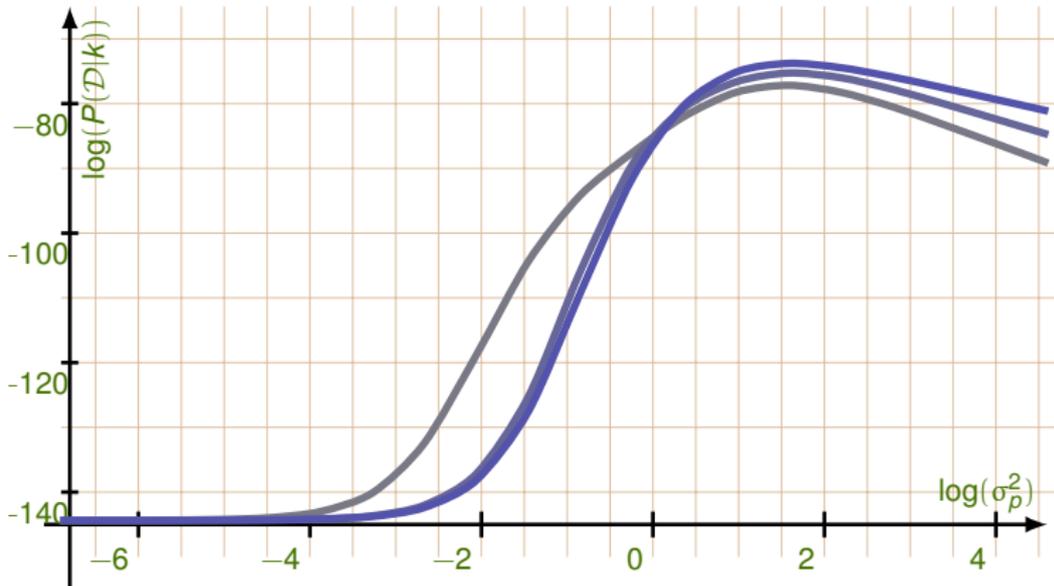
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



Aim: test different models.

Polynomial families: $k = 5$ $k = 4$ $k = 3$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

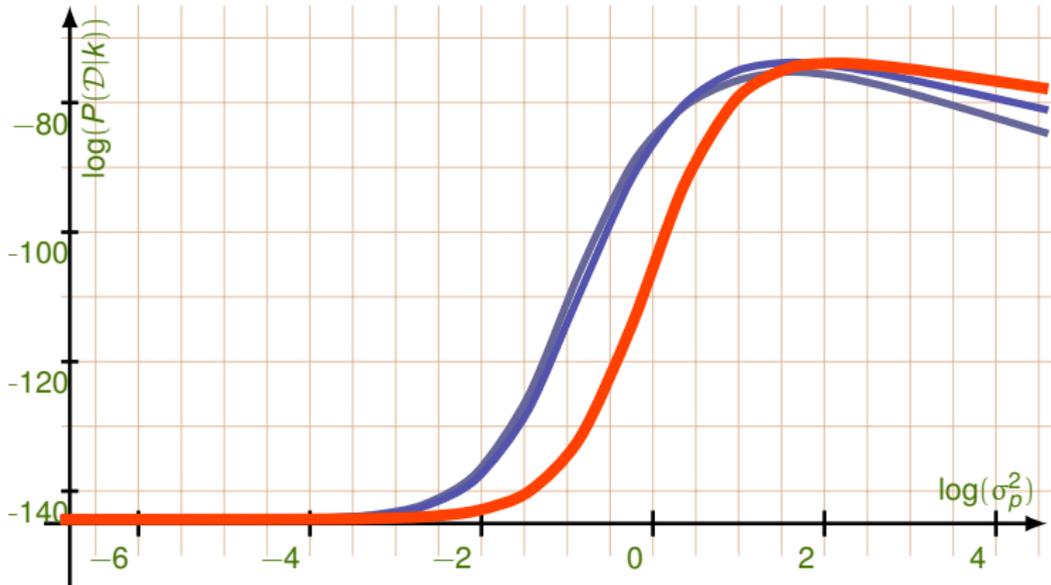
- Motivation
- Machine Learning
- Latent variable models

Estimation

- Maximum Likelihood
- Maximum a-posteriori
- Bayesian Estimation

Unsupervised

- General concepts
- Principal Components
- Independent Components
- Mixture Models



Aim: test different models.

Polynomial families: $k = 4$ $k = 3$ $k = 2$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

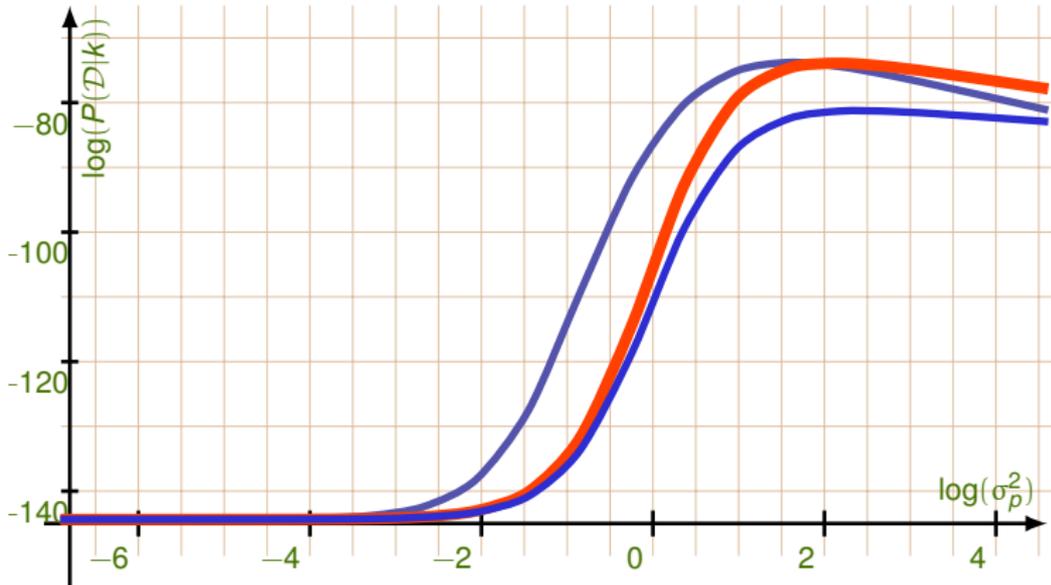
- Motivation
- Machine Learning
- Latent variable models

Estimation

- Maximum Likelihood
- Maximum a-posteriori
- Bayesian Estimation

Unsupervised

- General concepts
- Principal Components
- Independent Components
- Mixture Models



Aim: test different models.

Polynomial families: $k = 3$ $k = 2$ $k = 1$.



- M.L. and M.A.P. estimates provide **single** solutions.
- Point estimates lack the assessment of un/certainty.
- Better solution:
for a query \mathbf{x}_* , the system output is **probabilistic**:

$$\mathbf{x}_* \Rightarrow p(y_* | \mathbf{x}_*, \mathcal{F})$$

- Tool:
go beyond the M.A.P. solution and use the **a-posteriori distribution** of the parameters.



We again use Bayes' rule:

$$p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{F}) = \frac{P(\mathcal{D}|\boldsymbol{\theta})p_0(\boldsymbol{\theta})}{p(\mathcal{D}|\mathcal{F})} \quad \text{with } p(\mathcal{D}|\mathcal{F}) = \int_{\Omega} d\boldsymbol{\theta} P(\mathcal{D}|\boldsymbol{\theta})p_0(\boldsymbol{\theta}).$$

and **exploit** the whole posterior distribution of the parameters.

A-posteriori parameter estimates

We operate with $p_{\text{post}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{F})$ and use the total probability rule:

$$p(y_*|\mathcal{D}, \mathcal{F}) = \sum_{\boldsymbol{\theta}_\ell \in \Omega_{\boldsymbol{\theta}}} p(y_*|\boldsymbol{\theta}_\ell, \mathcal{F}) p_{\text{post}}(\boldsymbol{\theta}_\ell)$$

in assessing system output.



Given the data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ estimate the linear fit:

$$y = \theta_0 + \sum_{i=1}^d \theta_i x_i = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}^T \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \stackrel{\text{def}}{=} \boldsymbol{\theta}^T \mathbf{x}$$

Gaussian distributions noise and prior:

$$\epsilon = y_n - \boldsymbol{\theta}^T \mathbf{x}_n \sim \mathcal{N}(0, \sigma_n^2)$$
$$\boldsymbol{w} \sim \mathcal{N}(0, \boldsymbol{\Sigma}_0)$$



Goal: compute the posterior distribution $p_{\text{post}}(\boldsymbol{\theta})$.

$$p_{\text{post}}(\boldsymbol{\theta}) \propto p_0(\boldsymbol{\theta}) p(\mathcal{D}|\boldsymbol{\theta}, \mathcal{F}) = p_0(\boldsymbol{\theta}|\boldsymbol{\Sigma}_0) \prod_{n=1}^N P(y_n|\boldsymbol{\theta}^T \mathbf{x}_n)$$

$$\begin{aligned} -2 \log(p_{\text{post}}(\boldsymbol{\theta})) &= K_{\text{post}} + \frac{1}{\sigma_n^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) + \boldsymbol{\theta}^T \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\theta} \\ &= \boldsymbol{\theta}^T \left(\frac{1}{\sigma_n^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right) \boldsymbol{\theta} - \frac{2}{\sigma_n^2} \boldsymbol{\theta}^T \mathbf{X}^T \mathbf{y} + K'_{\text{post}} \\ &= (\boldsymbol{\theta} - \boldsymbol{\mu}_{\text{post}})^T \boldsymbol{\Sigma}_{\text{post}}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_{\text{post}}) + K''_{\text{post}} \end{aligned}$$

and by identification

$$\boldsymbol{\Sigma}_{\text{post}} = \left(\frac{1}{\sigma_n^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1} \quad \text{and} \quad \boldsymbol{\mu}_{\text{post}} = \boldsymbol{\Sigma}_{\text{post}} \frac{\mathbf{X}^T \mathbf{y}}{\sigma_n^2}$$



Bayesian linear model

The posterior distribution for the parameters is a Gaussian with parameters

$$\Sigma_{\text{post}} = \left(\frac{1}{\sigma_n^2} \mathbf{X}^T \mathbf{X} + \Sigma_0^{-1} \right)^{-1} \quad \text{and} \quad \mu_{\text{post}} = \Sigma_{\text{post}} \frac{\mathbf{X}^T \mathbf{y}}{\sigma_n^2}$$

Point estimates from keeping :

- M.L. if we take $\Sigma_0 \rightarrow \infty$ and considering only μ_{post} .
- M.A.P if we approximate the distribution with a single value at the maximum, *i.e.* μ_{post} .

**Prediction** for new values \mathbf{x}_* :

- use the likelihood $P(y_*|\mathbf{x}_*, \boldsymbol{\theta}, \mathcal{F})$,
- and the **posterior** for $\boldsymbol{\theta}$
- and Bayes' rule.

The steps:

$$\begin{aligned} p(y_*|\mathbf{x}_*, \mathcal{D}, \mathcal{F}) &= \int_{\Omega_{\boldsymbol{\theta}}} d\boldsymbol{\theta} p(y_*|\mathbf{x}_*, \boldsymbol{\theta}, \mathcal{F}) p_{\text{post}}(\boldsymbol{\theta} | \mathcal{D}, \mathcal{F}) \\ &= \int_{\Omega_{\boldsymbol{\theta}}} d\boldsymbol{\theta} \exp \left[-\frac{1}{2} \left(K_* + \frac{(y_* - \boldsymbol{\theta}^T \mathbf{x}_*)^2}{\sigma_n^2} + (\boldsymbol{\theta} - \boldsymbol{\mu}_{\text{post}})^T \boldsymbol{\Sigma}_{\text{post}}^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_{\text{post}}) \right) \right] \\ &= \int_{\Omega_{\boldsymbol{\theta}}} d\boldsymbol{\theta} \exp \left[-\frac{1}{2} \left(K_* + \frac{y_*^2}{\sigma_n^2} - \mathbf{a}^T \mathbf{C}^{-1} \mathbf{a} + Q(\boldsymbol{\theta}) \right) \right] \end{aligned}$$

where

$$\mathbf{a} = \frac{\mathbf{x}_* y_*}{\sigma_n^2} + \boldsymbol{\Sigma}_{\text{post}}^{-1} \boldsymbol{\mu}_{\text{post}} \quad \mathbf{C} = \frac{\mathbf{x}_* \mathbf{x}_*^T}{\sigma_n^2} + \boldsymbol{\Sigma}_{\text{post}}$$



Integrating out the quadratic in θ :

Predictive distribution at \mathbf{x}_*

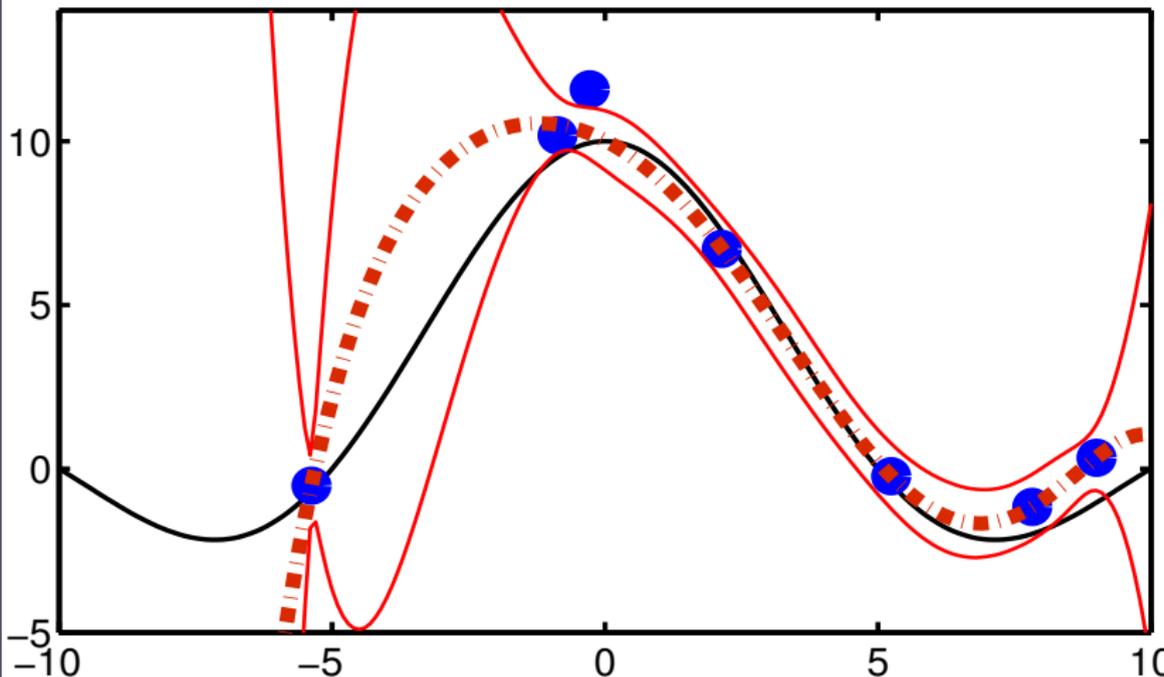
$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathcal{D}, \mathcal{F}) &= \exp \left[-\frac{1}{2} \left(K_* + \frac{(y_* - \mathbf{x}_*^T \boldsymbol{\mu}_{\text{post}})^2}{\sigma_n^2 + \mathbf{x}_*^T \boldsymbol{\Sigma}_{\text{post}}^{-1} \mathbf{x}_*} \right) \right] \\ &= \mathcal{N} \left(y_* \mid \mathbf{x}_*^T \boldsymbol{\mu}_{\text{post}}, \sigma_n^2 + \mathbf{x}_*^T \boldsymbol{\Sigma}_{\text{post}} \mathbf{x}_* \right) \end{aligned}$$

With the predictive distribution we:

- measure the variance of the prediction for each point:
 $\sigma_*^2 = \sigma_n^2 + \mathbf{x}_*^T \boldsymbol{\Sigma}_{\text{post}} \mathbf{x}_*$;
- sample from the parameters and plot the candidate predictors.



Pol. 6 – $N.\text{var}\sigma^2 = 1$



The errors are the symmetric thin lines.



Bayesian example

Predictive samples

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

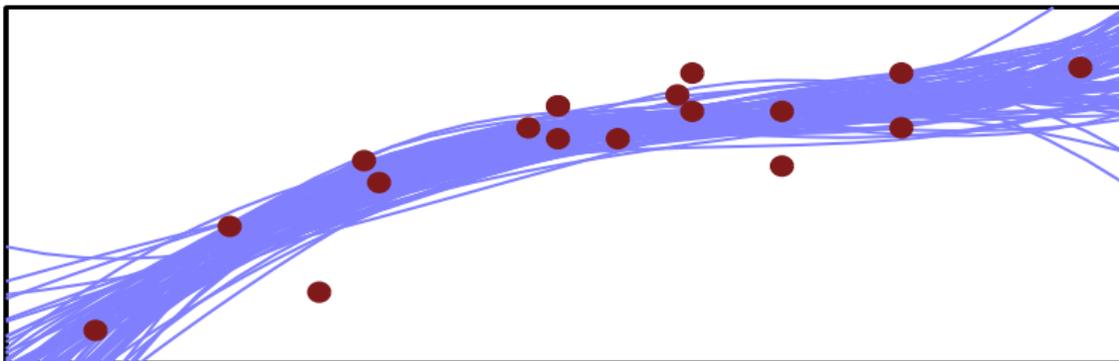
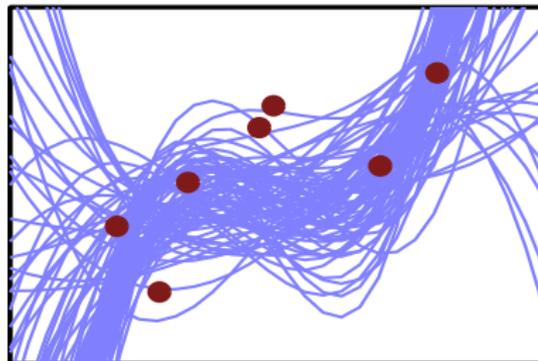
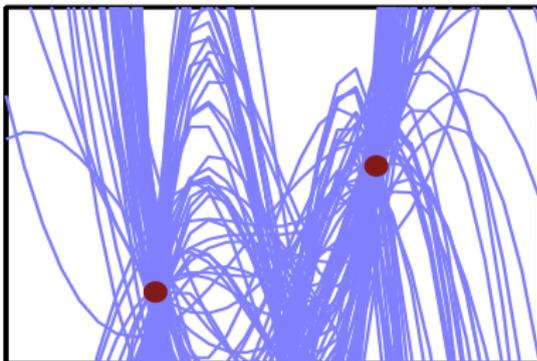
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



Third order polynomials are used to approximate the data.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

When computing $p_{\text{post}}(\boldsymbol{\theta}|\mathcal{D}, \mathcal{F})$ we assumed that the posterior **can be represented** analytically.

This is not the case.

Approximations are needed for the

- posterior distribution
- predictive distribution

In Bayesian modelling an important issue is **how** we approximate the posterior distribution.



Complete specification of the model

Can include prior beliefs about the model.

Accurate predictions

Can compute the posterior probabilities for each test location.

Computational cost

Using models for prediction can be difficult and expensive in time and memory.



Complete specification of the model

Can include prior beliefs about the model.

Accurate predictions

Can compute the posterior probabilities for each test location.

Computational cost

Using models for prediction can be difficult and expensive in time and memory.

Bayesian models

Flexible and accurate – **if** priors about the model are used.



Outline

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- 1 Modelling Data
- 2 Estimation methods
- 3 Unsupervised Methods**
 - General concepts
 - Principal Components
 - Independent Components
 - Mixture Models



Unsupervised setting

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

- Data can be unlabeled, *i.e.* no values y are associated to an input \mathbf{x} .
- We want to “extract” information from $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.
- We assume that the data – although high-dimensional – span a **much smaller dimensional manifold**.
- Task is to find the subspace corresponding to the *data span*.



Models in unsupervised learning

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

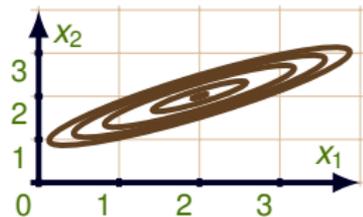
Principal Components

Independent Components

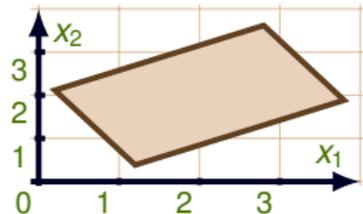
Mixture Models

It is again important the **model of the data**:

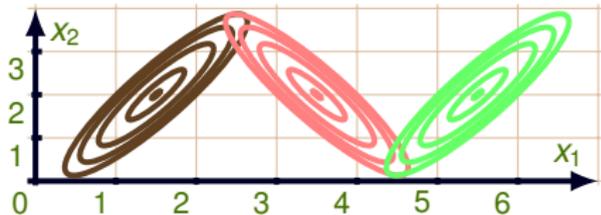
- Principal Components;



- Independent Components;



- Mixture models;





The PCA model

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

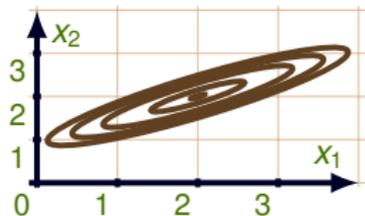
General concepts

Principal Components

Independent Components

Mixture Models

- Simple data structure.
- Spherical cluster that is:
 - translated;
 - scaled;
 - rotated.



We aim to find the principal directions of the data spread.

Principal direction:

the direction u along which the data preserves most of its *variance*.



Principal direction:

$$\mathbf{u} = \operatorname{argmax}_{\|\mathbf{u}\|=1} \frac{1}{2N} \sum_{n=1}^N (\mathbf{u}^T \mathbf{x}_n - \mathbf{u} \bar{\mathbf{x}})^2$$

we pre-process: $\bar{\mathbf{x}} = \mathbf{0}$. Replacing the empirical covariance with $\Sigma_{\mathbf{x}}$:

$$\begin{aligned} \mathbf{u} &= \operatorname{argmax}_{\|\mathbf{u}\|=1} \frac{1}{2N} \sum_{n=1}^N (\mathbf{u}^T \mathbf{x}_n - \mathbf{u} \bar{\mathbf{x}})^2 \\ &= \operatorname{argmax}_{\mathbf{u}, \lambda} \frac{1}{2} \mathbf{u}^T \Sigma_{\mathbf{x}} \mathbf{u} - \lambda (\|\mathbf{u}\|^2 - 1) \end{aligned}$$

with λ the Lagrange multiplier. Differentiating w.r.t \mathbf{u} :

$$\Sigma_{\mathbf{x}} \mathbf{u} - \lambda \mathbf{u} = \mathbf{0}$$



The optimum solution **must obey**:

$$\Sigma_x \mathbf{u} = \lambda \mathbf{u}$$

The **eigendecomposition** of the covariance matrix.

$(\lambda_*, \mathbf{u}_*)$ is an eigenvalue, eigenvector of the system.

If we replace back, the value of the expression is λ_* .



Optimal solution when $\lambda_* = \lambda_{max}$.

Principal direction:

The eigenvector \mathbf{u}_{max} corresponding to the largest eigenvalue of the system.



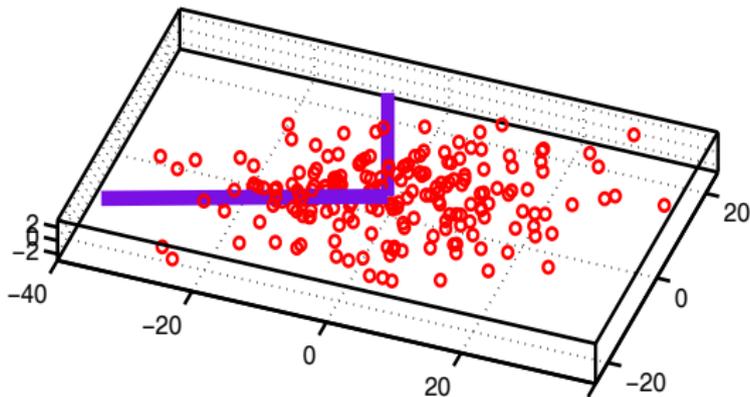
How is this used in data mining?

Assume that data is:

- jointly Gaussian:

$$\mathbf{x} = \mathbf{N}(\mathbf{m}_x, \Sigma_x),$$

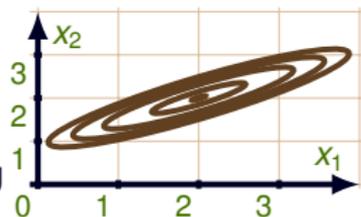
- high-dimensional;
- only few (2) directions are relevant.





How is this used in data mining?

- Subtracting mean.
- Eigendecomposition.
- Selecting the K eigenvectors corresponding to the K largest values.
- Computing the K projections: $z_{nl} = \mathbf{x}_n^T \mathbf{u}_l$.



The projection using matrix $\mathbf{P} \stackrel{\text{def}}{=} [\mathbf{u}_1, \dots, \mathbf{u}_K]^T$:

$$\mathbf{Z} = \mathbf{XP}$$

and \mathbf{z}_n can be used as a compact representation of \mathbf{x}_n .



Reconstruction:

$$\mathbf{x}'_n = \sum_{\ell=1}^K z_{n\ell} \mathbf{u}_\ell \quad \text{or, with matrix notation: } \mathbf{X}' = \mathbf{Z}\mathbf{P}^T$$

PCA projection analysis:

$$\begin{aligned} E_{PCA} &= \frac{1}{N^2} \sum_{n=1}^N (\mathbf{x}_n - \mathbf{x}'_n)^2 = \frac{1}{N^2} \text{tr} [(\mathbf{X} - \mathbf{X}')^T (\mathbf{X} - \mathbf{X}')] \\ &= \text{tr} [\boldsymbol{\Sigma}_x - \mathbf{P}^T \boldsymbol{\Sigma}_z \mathbf{P}] \\ &= \text{tr} [\mathbf{U} (\text{diag}(\lambda_1, \dots, \lambda_d) - \text{diag}(\lambda_1, \dots, \lambda_K, 0, \dots)) \mathbf{U}^T] \\ &= \text{tr} [\mathbf{U}^T \mathbf{U} \text{diag}(0, \dots, 0, \lambda_{K+1}, \dots, \lambda_d)] \\ &= \sum_{\ell=1}^{d-K} \lambda_{K+\ell} \end{aligned}$$



PCA reconstruction error:

The error made using the PCA directions:

$$E_{PCA} = \sum_{\ell=K+1}^{d-K} \lambda_{K+\ell}$$

PCA properties:

- PCA system orthonormal: $\mathbf{u}_\ell^T \mathbf{u}_r = \delta_{\ell-r}$
- Reconstruction **fast**.
- Spherical assumption critical.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

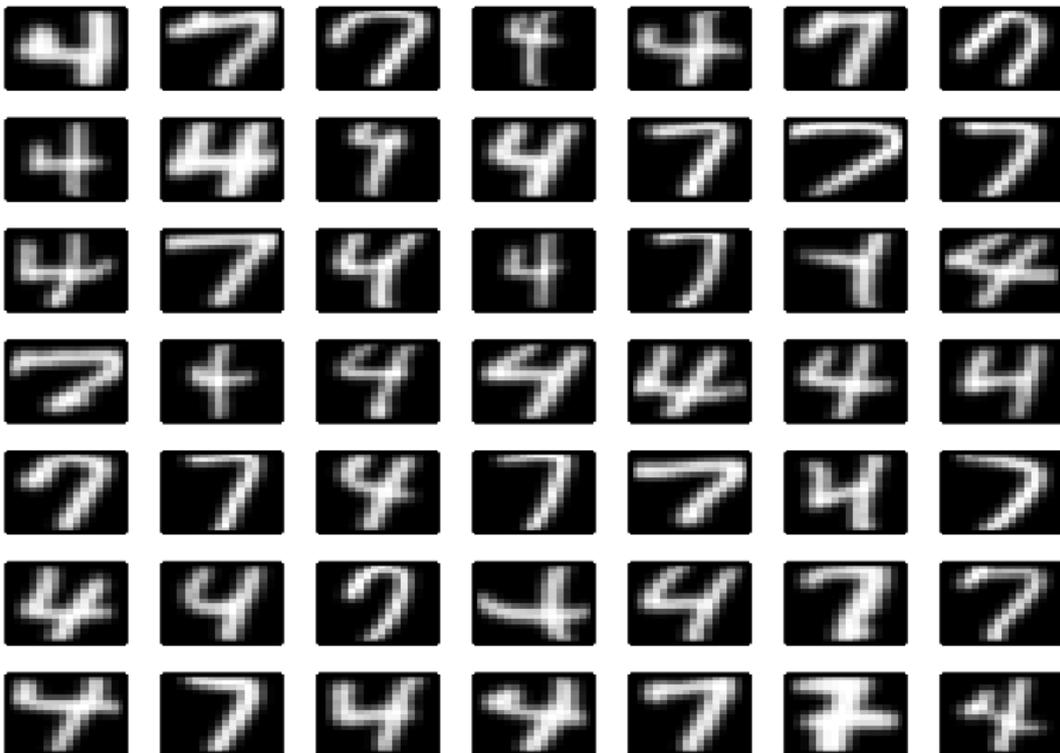
General concepts

Principal Components

Independent Components

Mixture Models

USPS digits – testbed for several models.



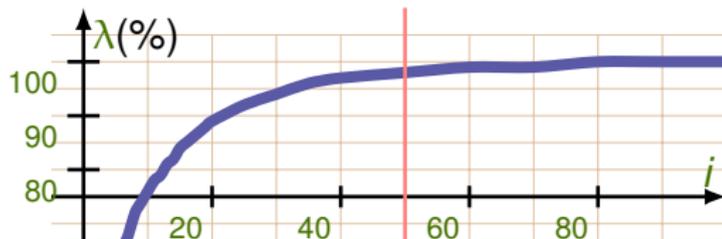


USPS characteristics:

- handwritten data centered and scaled;
- ≈ 10.000 items of 16×16 grayscale images;

We plot

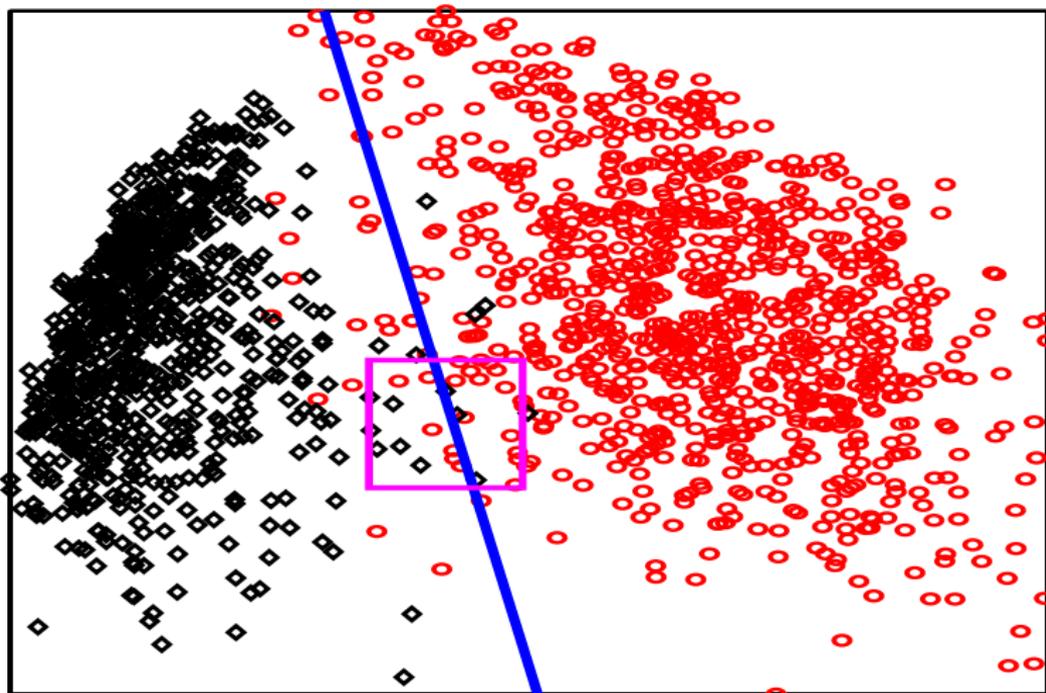
$$k_r = \sum_{\ell=1}^r \lambda_{\ell}$$



Conclusion for the USPS set:

- The normalised $\lambda_1 = 0.24 \Rightarrow \mathbf{u}_1$ accounts for 24% of the data.
- at ≈ 10 more than 70% of variance is explained.
- at ≈ 50 more than 98%
 \Rightarrow 50 numbers instead of 256.

Visualisation application:



Visualisation along the **first two** eigendirections.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

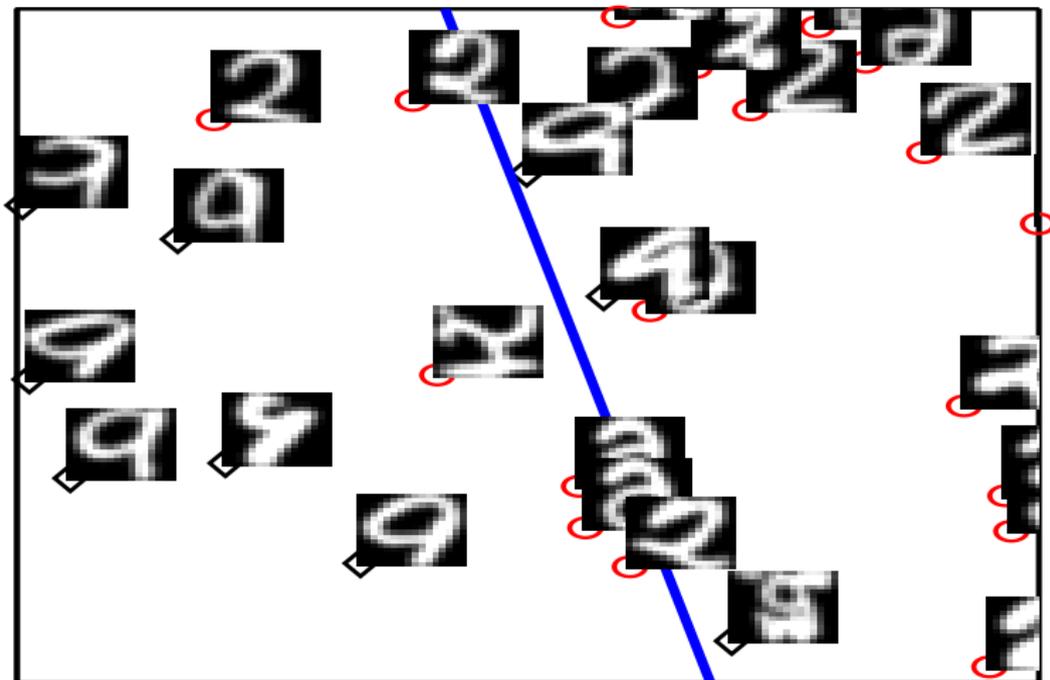
General concepts

Principal Components

Independent Components

Mixture Models

Visualisation application:



Detail.



The ICA model

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

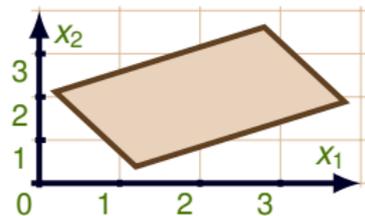
Independent Components

Mixture Models

Start from the PCA:

$$\mathbf{x} = \mathbf{P}\mathbf{z}$$

is a **generative model** for the data.



We assumed that

- \mathbf{z} i.i.d. Gaussian random variables $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \text{diag}(\lambda_\ell))$;
- $\Rightarrow \mathbf{x}$ are not independent;
- $\Rightarrow \mathbf{z}$ are **Gaussian** sources;

In most of real data:

- Sources are **not Gaussian**.
- But sources are **independent**.
- **We exploit that!**



The following model assumption:

$$x = As$$

where

- z independent **sources**;
- A linear **mixing matrix**;

Looking for matrix B that recovers the sources:

$$s' \stackrel{\text{def}}{=} Bx = B(As) = (BA)s$$

i.e. (BA) is unity up to a permutation and scaling
but retains **independence**.



In practice:

$$\mathbf{s}' \stackrel{\text{def}}{=} \mathbf{B}\mathbf{x}$$

with $\mathbf{s} = [s_1, \dots, s_K]$ all independent sources.

Independence test: the KL-divergence between the **joint distribution** and the **marginals**

$$\mathbf{B} = \underset{\mathbf{B} \in \text{SO}_d}{\text{argmin}} \text{KL} (p(s_1, s_2) || p(s_1)p(s_2))$$

where SO_d is the group of matrices with $|\mathbf{B}| = 1$.

In ICA we are looking for matrix \mathbf{B} that minimises:

$$\sum_{\ell} \int_{\Omega_{\ell}} dp(s_{\ell}) \log p(s_{\ell}) - \int_{\Omega_{\ell}} dp(\mathbf{s}) \log(p(s_1, \dots, s_d))$$



Kullback-Leibler divergence

$$\text{KL}(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- is zero only and only if $p = q$,
- is *not* a measure of distance (but cloooooose to it!),
- Efficient when exponential families are used.

Short proof:

$$\begin{aligned} 0 &= \log 1 = \log \left(\sum_x q(x) \right) = \log \left(\sum_x p(x) \frac{q(x)}{p(x)} \right) \\ &\geq \sum_x p(x) \log \left(\frac{q(x)}{p(x)} \right) = -\text{KL}(p||q) \end{aligned}$$

$$\Rightarrow \text{KL}(p||q) \geq 0$$



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

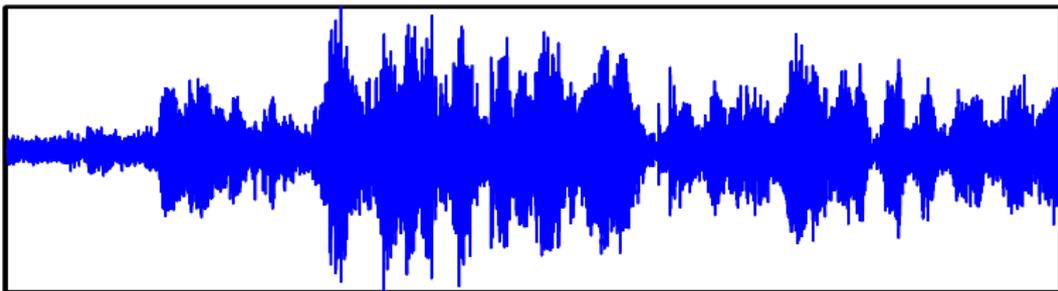
Principal Components

Independent Components

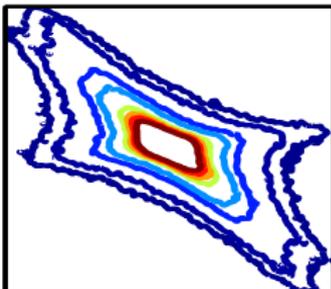
Mixture Models

Separation of source signals:

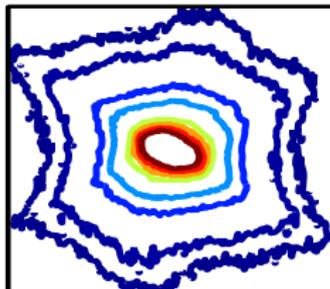
Mixture



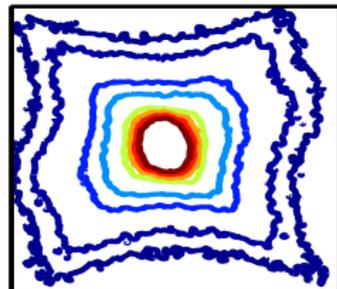
m2 m4



m1 m3



m3 m4





Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

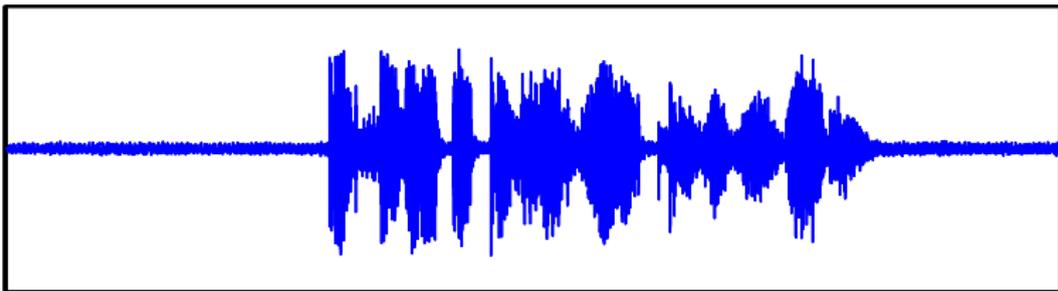
Principal Components

Independent Components

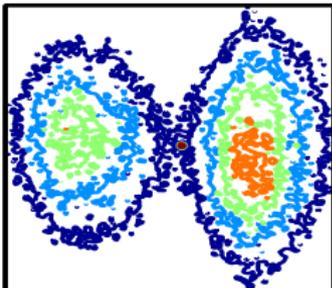
Mixture Models

Results of separation:

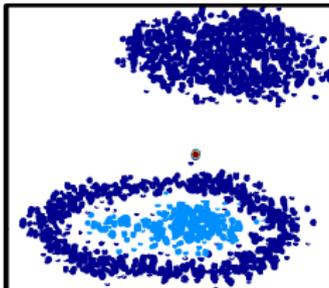
Source



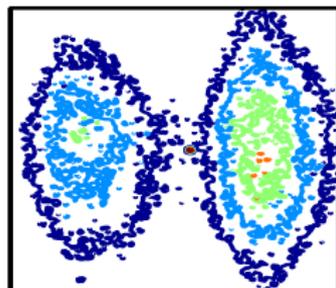
m2 m4



m1 m3



m3 m4



FastICA package



Applications of ICA

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

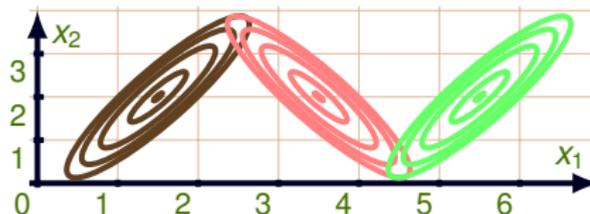
Mixture Models

Applications:

- **Coctail party problem;**
Separates noisy and multiple sources from multiple observations.
- **Fetus ECG;**
Separation of the ECG signal of a fetus from its mother's ECG.
- **MEG recordings;**
Separation of MEG "sources".
- **Financial data;**
Finding hidden factors in financial data.
- **Noise reduction;**
Noise reduction in natural images.
- **Interference removal;**
Interference removal from CDMA – Code-division multiple access – communication systems.



- The data structure is **more complex**.
- More than a single source for data.



The mixture model:

$$P(\mathbf{x}|\Sigma) = \sum_{k=1}^K \pi_k \rho_k(\mathbf{x}|\mu_k, \Sigma_k) \quad (1)$$

where:

π_1, \dots, π_K – mixing components.

$\rho_k(\mathbf{x}|\mu_k, \Sigma_k)$ – density of a component.

The components are usually called clusters.



The generation process reflects the assumptions about the model.

The data generation:

- first we select **from which** component,
- then we sample from the component's density function.

When modelling data we **do not know**:

- Which point belongs to which cluster.
- What are the parameters for each density function.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

Old Faithful geyser in the Yellowstone National park.

Characterised by:

- intense eruptions;
- differing times between them.

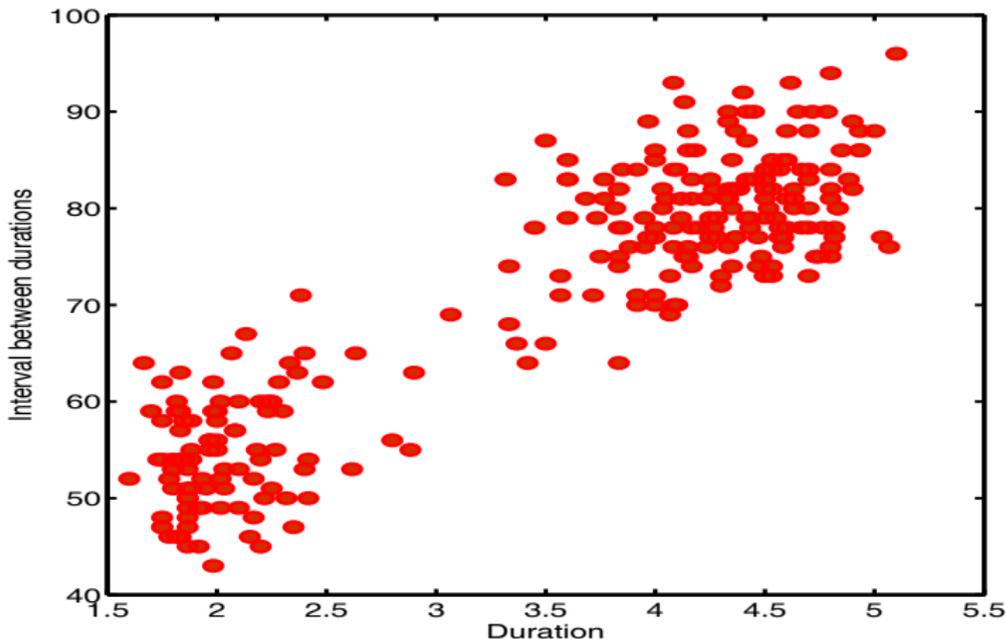
Rule:

Duration is 1.5 to 5 minutes.

The length of eruption helps determine the interval.

If an eruption lasts less than 2 minutes the interval will be around 55 minutes. If the eruption last 4.5 minutes the interval may be around 88 minutes.





- The longer the duration, the longer the interval.
- The linear relation $l = \theta_0 + \theta_1 d$ is not the best.
- There are only a very few eruptions lasting ≈ 3 minutes.



The mixture model

Assumptions:

- We know the family of individual density functions:
These density functions are parametrised with a few parameters.
- The densities are easily identifiable:
If we knew which data belongs to which cluster, the density function is easily identifiable.

Gaussian densities are often used – fulfill both “conditions”.



The mixture model

II

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

The Gaussian mixture model:

$$p(\mathbf{x}) = \pi_1 N_1(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \pi_2 N_2(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

for **known** densities
(centres and ellipses):

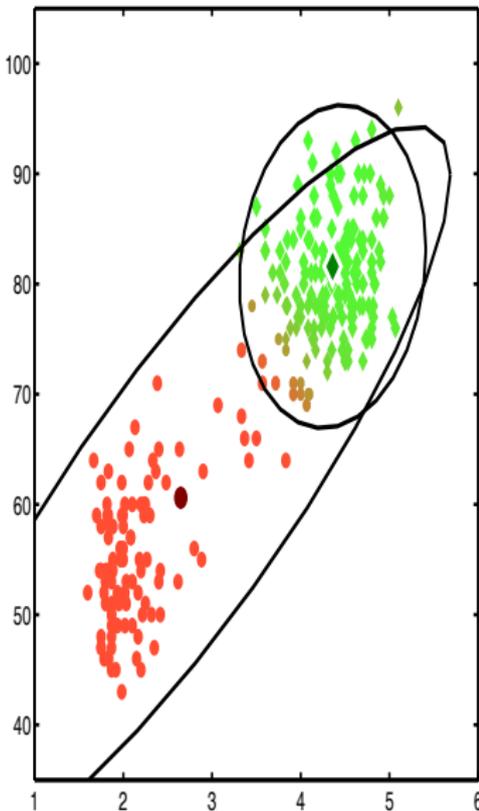
$$p(\mathbf{x}_n|k) = \frac{N_k(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) p(k)}{\sum_{\ell} N_{\ell}(\mathbf{x}_n|\boldsymbol{\mu}_{\ell}, \boldsymbol{\Sigma}_{\ell}) p(\ell)}$$

i.e. we know the **probability** that data comes from cluster k
(shades from red to green).

For \mathcal{D} :

\mathbf{x}	$p(\mathbf{x} 1)$	$p(\mathbf{x} 2)$
\mathbf{x}_1	γ_{11}	γ_{12}
\vdots	\vdots	\vdots
\mathbf{x}_N	γ_{N1}	γ_{N2}

γ_{nl} – responsibility of \mathbf{x}_n in cluster ℓ .





The mixture model



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

When γ_{nl} known, the parameters are computed using the data weighted by their responsibilities:

$$(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \underset{\boldsymbol{\mu}, \boldsymbol{\Sigma}}{\operatorname{argmax}} \prod_{n=1}^N (\mathcal{N}_k(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}))^{\gamma_{nk}}$$

for all k .

This means:

$$(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \leftarrow \sum_n \gamma_{nk} \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

When making inference

Have to find the responsibility vector **and** the parameters of the mixture.



The mixture model



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

When γ_{ne} known, the parameters are computed using the data weighted by their responsibilities:

$$(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \operatorname{argmax}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} \prod_{n=1}^N (\mathbb{N}_k(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}))^{\gamma_{nk}}$$

for all k .

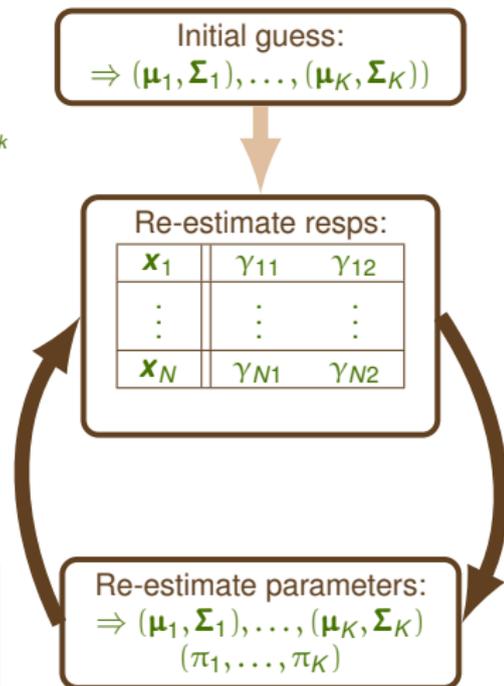
This means:

$$(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \leftarrow \sum_n \gamma_{nk} \log \mathbb{N}(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

When making inference

Have to find the responsibility vector **and** the parameters of the mixture.

Given data \mathcal{D} :





Responsibilities γ

The additional **latent** variables needed to help computation.

In the mixture model:

- goal is to *fit* model to data;
- which submodel gets a particular data;

Achieved by the maximisation of the log-likelihood function.



The EM algorithm

Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models

$$(\boldsymbol{\pi}, \boldsymbol{\Theta}) = \operatorname{argmax} \sum_n \log \left[\sum_{\ell} \pi_{\ell} \mathcal{N}_{\ell}(\mathbf{x}_n | \boldsymbol{\mu}_{\ell}, \boldsymbol{\Sigma}_{\ell}) \right]$$

$\boldsymbol{\Theta} = [\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K]$ is the vector of parameters;
 $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$ the shares of the factors;

Problem with optimisation:

The parameters are not separable due to the sum within the logarithm.

Solution:

Use an approximation.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation
Machine Learning
Latent variable models

Estimation

Maximum Likelihood
Maximum a-posteriori
Bayesian Estimation

Unsupervised

General concepts
Principal Components
Independent Components
Mixture Models

$$\begin{aligned}\log P(\mathcal{D}|\boldsymbol{\pi}, \boldsymbol{\Theta}) &= \sum_n \log \left[\sum_{\ell} \pi_{\ell} N_{\ell}(\mathbf{x}_n | \boldsymbol{\mu}_{\ell}, \boldsymbol{\Sigma}_{\ell}) \right] \\ &= \sum_n \log \left[\sum_{\ell} p_{\ell}(\mathbf{x}_n, \ell) \right]\end{aligned}$$

Use Jensen's inequality:

$$\begin{aligned}\log \left(\sum_{\ell} p_{\ell}(\mathbf{x}_n, \ell | \theta_{\ell}) \right) &= \log \left(\sum_{\ell} q_n(\ell) \frac{p_{\ell}(\mathbf{x}_n, \ell | \theta_{\ell})}{q_n(\ell)} \right) \\ &\geq \sum_{\ell} q_n(\ell) \log \left(\frac{p_{\ell}(\mathbf{x}_n, \ell)}{q_n(\ell)} \right)\end{aligned}$$

for **any** $[q_n(1), \dots, q_n(\ell)]$.



Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

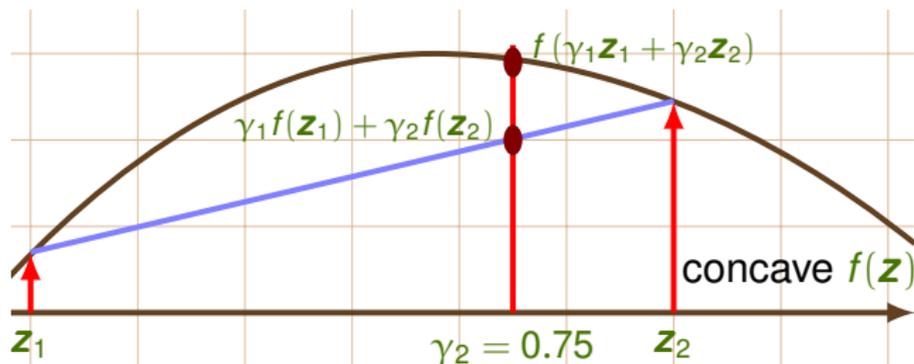
Unsupervised

General concepts

Principal Components

Independent Components

Mixture Models



Jensen's Inequality

For any concave $f(z)$, any z_1 and z_2 , and any $\gamma_1, \gamma_2 > 0$ such that $\gamma_1 + \gamma_2 = 1$:

$$f(\gamma_1 z_1 + \gamma_2 z_2) \geq \gamma_1 f(z_1) + \gamma_2 f(z_2)$$



$$\log \left(\sum_{\ell} p_{\ell}(\mathbf{x}_n, \ell | \theta_{\ell}) \right) \geq \sum_{\ell} q_n(\ell) \log \left(\frac{p_{\ell}(\mathbf{x}_n, \ell)}{q_n(\ell)} \right)$$

for any distribution $q_n(\cdot)$.

Replacing with the right-hand side, we have:

$$\begin{aligned} \log P(\mathcal{D} | \boldsymbol{\pi}, \boldsymbol{\Theta}) &\geq \sum_n \sum_{\ell} q_n(\ell) \log \frac{p_{\ell}(\mathbf{x}_n | \boldsymbol{\theta}_{\ell})}{q_n(\ell)} \\ &\geq \sum_{\ell} \left[\sum_n q_n(\ell) \log \frac{p_{\ell}(\mathbf{x}_n | \boldsymbol{\theta}_{\ell})}{q_n(\ell)} \right] = \mathcal{L} \end{aligned}$$

and therefore the optimisation w.r.to cluster parameters separate.

$$\partial_{\ell} \Rightarrow 0 = \sum_n q_n(\ell) \frac{\partial \log p_{\ell}(\mathbf{x}_n | \boldsymbol{\theta}_{\ell})}{\partial \boldsymbol{\theta}_{\ell}}$$

For distributions from exponential family optimisation is easy.



- **any** set of distributions $q_1(\ell), \dots, q_N(\ell)$ provides a lower bound to the log-likelihood.
- We should choose the distributions so that they are the closest to the **current** parameter set.

We assume the **parameters** have the value θ_0 .

Want to minimise the difference:

$$\log P(\mathbf{x}_n, \ell | \theta_\ell^0) - \mathcal{L} = \sum_{\ell} q_n(\ell) \log P(\mathbf{x}_n, \ell | \theta_\ell^0) - \sum_{\ell} q_n(\ell) \log \frac{p_{\ell}(\mathbf{x}_n, \ell | \theta_{\ell}^0)}{q_n(\ell)}$$
$$\sum_{\ell} q_n(\ell) \log \frac{P(\mathbf{x}_n, \ell | \theta_{\ell}^0) q_n(\ell)}{p_{\ell}(\mathbf{x}_n | \theta_{\ell}^0)}$$

and observe that by **setting**

$$q_n(\ell) = \frac{p_{\ell}(\mathbf{x}_n | \theta_{\ell}^0)}{P(\mathbf{x}_n, \ell | \theta_{\ell}^0)}$$

we have $\sum_{\ell} q_n(\ell) \theta = 0$.



The EM algorithm:

Init – initialise model parameters;

E step – compute the responsibilities $\gamma_{nl} = q_n(\ell)$;

M step – for each k optimize

$$0 = \sum_n q_n(\ell) \frac{\partial \log p_\ell(\mathbf{x}_n | \boldsymbol{\theta}_\ell)}{\partial \boldsymbol{\theta}_\ell}$$

repeat – goto the **E step**.



EM application

Probabilistic Data Mining

Lehel Csató

Modelling Data

- Motivation
- Machine Learning
- Latent variable models

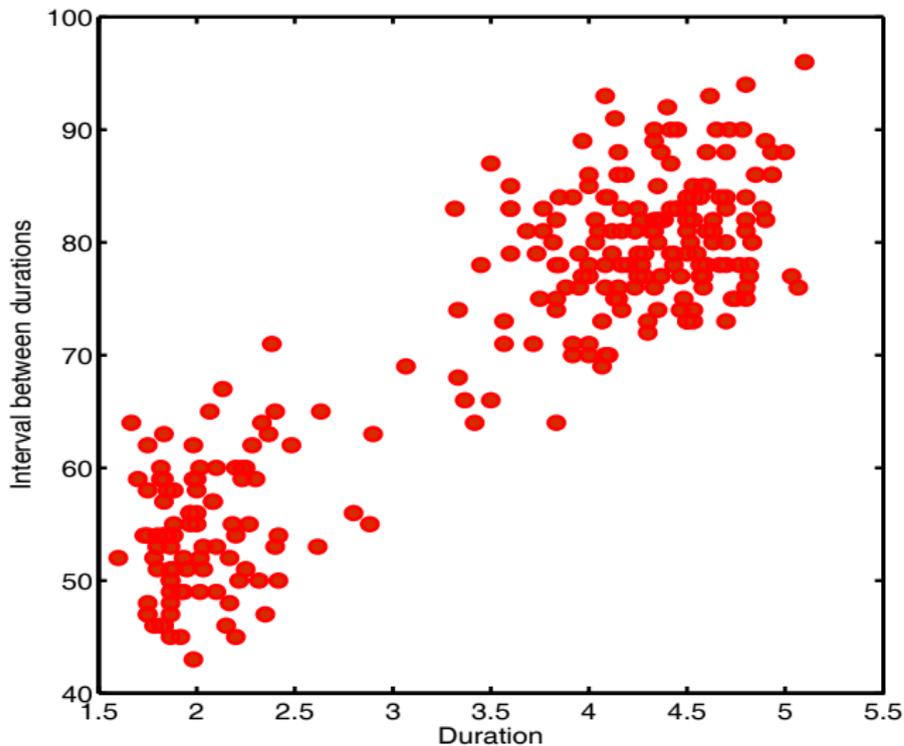
Estimation

- Maximum Likelihood
- Maximum a-posteriori
- Bayesian Estimation

Unsupervised

- General concepts
- Principal Components
- Independent Components
- Mixture Models

Old faithful:





Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

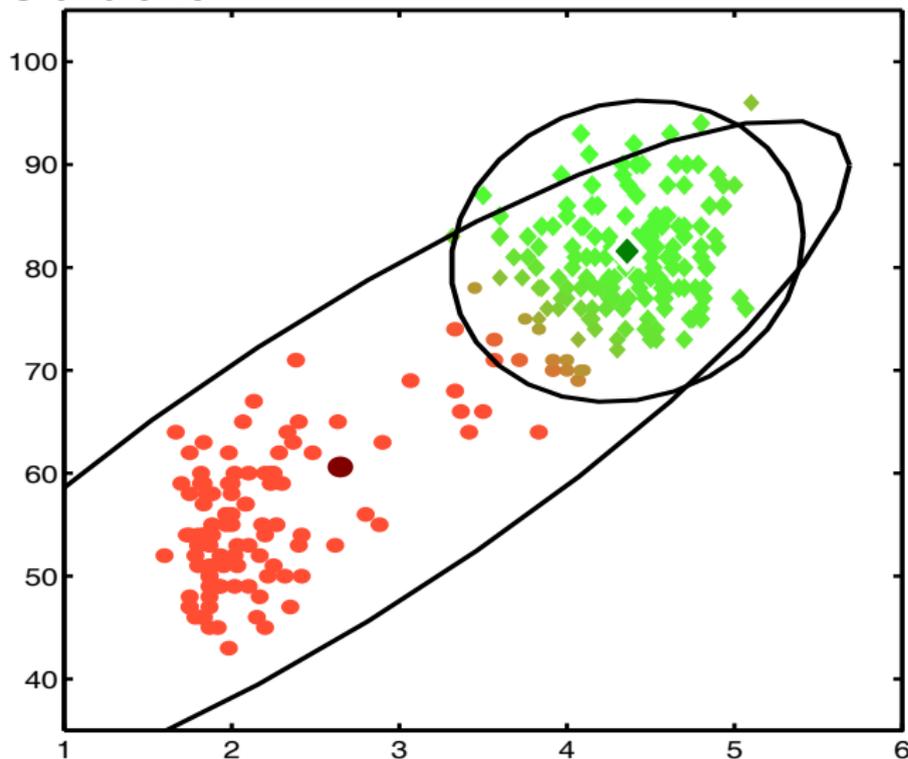
General concepts

Principal Components

Independent Components

Mixture Models

Old faithful:





Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

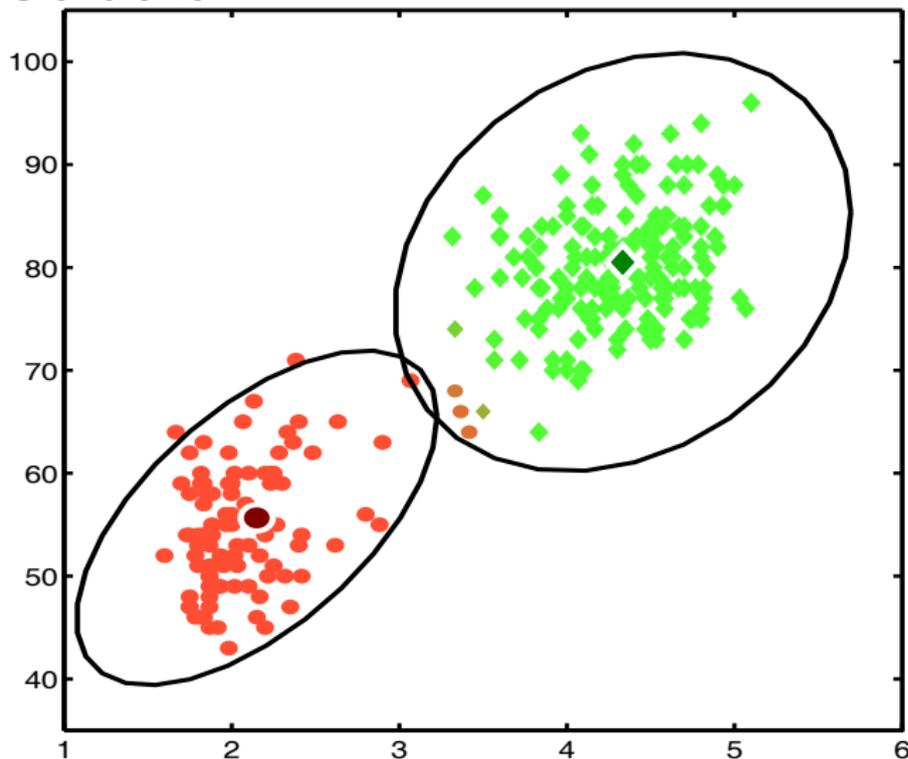
General concepts

Principal Components

Independent Components

Mixture Models

Old faithful:





Probabilistic Data Mining

Lehel Csató

Modelling Data

Motivation

Machine Learning

Latent variable models

Estimation

Maximum Likelihood

Maximum a-posteriori

Bayesian Estimation

Unsupervised

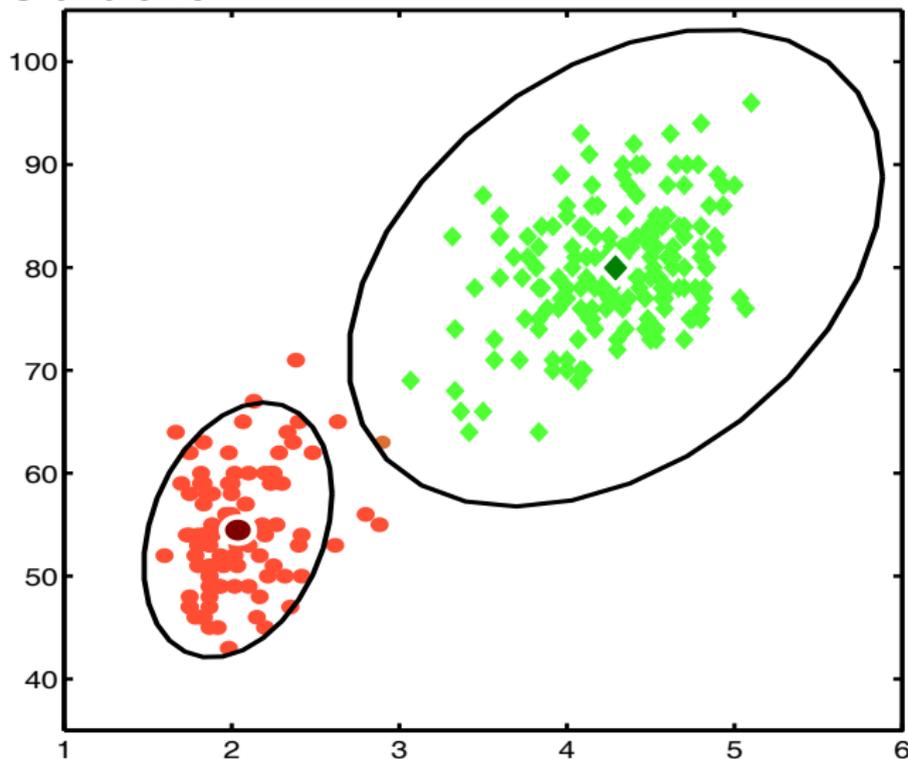
General concepts

Principal Components

Independent Components

Mixture Models

Old faithful:





References

Probabilistic Data
Mining

Lehel Csató

References



J. M. Bernardo and A. F. Smith.
Bayesian Theory.
John Wiley & Sons, 1994.



C. M. Bishop.
Pattern Recognition and Machine Learning.
Springer Verlag, New York, N.Y., 2006.



T. M. Cover and J. A. Thomas.
Elements of Information Theory.
John Wiley & Sons, 1991.



A. P. Dempster, N. M. Laird, and D. B. Rubin.
Maximum likelihood from incomplete data via the EM algorithm.
Journal of the Royal Statistical Society series B, 39:1–38, 1977.



T. Hastie, R. Tibshirani, és J. Friedman.
The Elements of Statistical Learning: Data Mining, Inference, and Prediction.
Springer Verlag, 2001.