

Inverse problems in Finance and Human sciences

Ali Mohammad-Djafari

Laboratoire des Signaux et Systèmes (L2S)

UMR8506 CNRS-CentraleSupélec-UNIV PARIS SUD

SUPELEC, 91192 Gif-sur-Yvette, France

<http://lss.centralesupelec.fr>

Email: djafari@lss.supelec.fr

<http://djafari.free.fr>

<http://publicationslist.org/djafari>

Workshop Inverse problems in Finance and Human sciences, ATU,
Tehran, Iran, September 24-25, 2016

Contents

1. Examples of inverse problems
 - ▶ Low dimensional case
 - ▶ High dimensional case
2. Basics of Bayesian inference
3. Bayes for Inverse Problems and Machine Learning
(Estimation, Prediction, Model Evaluation and selection)
4. Approximate Bayesian Computation (ABC)
 - ▶ Laplace approximation
 - ▶ Bayesian Information Criterion (BIC)
 - ▶ Variational Bayesian Approximation
 - ▶ Expectation Propagation (EP), Message Passing, MCMC, Exact Sampling, ...
5. Bayes for inverse problems
 - ▶ Traffic Management and Computed Tomography: A Linear problem
 - ▶ Differential Equations in Finance and Microwave imaging: A Bi-Linear or Non-Linear problem
6. Some canonical problems in Machine Learning
 - ▶ Regression, Classification and Model selection for classical and Big Data cases

Examples of inverse problems

1. Discrete cases examples
2. Continuous cases examples

Traffic management

$$\begin{array}{ccc} \textcolor{blue}{r}_1 \diamond & & \star \textcolor{blue}{c}_1 \\ \vdots & & \vdots \\ \sum_j \textcolor{red}{f}_{ij} = \textcolor{blue}{r}_i \diamond & \{\textcolor{red}{f}_{ij}\} & \star \textcolor{blue}{c}_j = \sum_i \textcolor{red}{f}_{ij} \\ \vdots & & \vdots \\ \textcolor{blue}{r}_I \diamond & & \star \textcolor{blue}{r}_J \end{array}$$

- ▶ M residential places each containing $\textcolor{blue}{r}_i, i = 1, \dots, I$ cars
- ▶ N working places each containing $\textcolor{blue}{c}_j, j = 1, \dots, J$ parking lots
- ▶ We want to estimate the numbers $\textcolor{red}{f}_{i,j}$ of cars going from the residential place i to working place j
- ▶ we know:

$$\sum_{j=1}^J \textcolor{red}{f}_{i,j} = \textcolor{blue}{r}_i, i = 1, \dots, I, \quad \sum_{i=1}^I \textcolor{red}{f}_{i,j} = \textcolor{blue}{c}_j, j = 1, \dots, J$$

- ▶ find $\textcolor{red}{f}_{i,j}$.

Traffic management: A very low dimensional and simple example

- $I = 2, J = 2$

$$\sum_{j=1}^2 f_{i,j} = r_i, i = 1, 2 \quad \sum_{i=1}^2 f_{i,j} = c_j, j = 1, 2$$

$$r_1 = 4, r_2 = 6, c_1 = 3, c_2 = 7$$

- Writing it differently: find $f_{i,j}$

$f_{1,1}$	$f_{1,2}$	4
$f_{2,1}$	$f_{2,2}$	6
3	7	

- A second example

$f_{1,1}$	$f_{1,2}$	$f_{1,3}$	9
$f_{2,1}$	$f_{2,2}$	$f_{2,3}$	10
3	7	11	

- Then, we extend it greater dimension

Traffic management

Think about:

- ▶ Has this problem a solution?
- ▶ Is the solution unique?
- ▶ How to find all the possible solutions?
- ▶ How we can select one of these solution?
- ▶ Minimum Norm solutions? $\min \sum_{i,j} f_{i,j}^2$ subject to the data constraints.
- ▶ What about minimizing:
 - ▶ l_1 norm: $\sum_{i,j} |f_{i,j}|$
 - ▶ l_α norm: $\sum_{i,j} |f_{i,j}|^\alpha$
 - ▶ Entropy: $-\sum_{i,j} f_{i,j} \log f_{i,j}$
- ▶ What if there are uncertainties in the data ?
- ▶ If you have to decide to construct for a fast road which one you propose?

Continuous equivalent problem

- Given two functions $r(y)$ and $c(x)$, find a function and $f(x, y)$ of 2 variables x and y such that:

$$\begin{cases} \int f(x, y) dx = r(y) \\ \int f(x, y) dy = c(x) \end{cases}$$

- When $r(y)$ and $c(x)$ are the marginal probability distribution of a joint probability distribution $f(x, y)$, we see the **Copula** theory

Prediction of gain

- ▶ Three days ago you gained 200 Euros, Two days ago you gained 100 Euros, Yesterday 200 Euros, Today 100 Euros.
- ▶ How much are you expecting to gain tomorrow? after tomorrow ?

t_i	-3	-2	-1	0	1	2
x_i	200	100	160	200	?	?

- ▶ Think about the following models:
 - ▶ $x(t_i) = \theta_0 + \theta_1(t_i) + \epsilon_i$
 - ▶ $x(t_i) = \theta_0 + \theta_1(t_i) + \theta_2(t_i^2) + \epsilon_i$
 - ▶ $x(t_i) = \theta_0 + \theta_1(t_i) + \theta_2(t_i^2) + \theta_3(t_i^3) + \epsilon_i$
 - ▶ ...
 - ▶ $x_i = \theta_1 x_{i-1} + \epsilon_i$
 - ▶ $x_i = \theta_1 x_{i-1} + \theta_2 x_{i-2} + \epsilon_i$
 - ▶ $x_i = \sum_k^K \theta_k x_{i-k} + \epsilon_i$
 - ▶ ...
 - ▶ $x_i = \theta_0 + \theta_1 \sin(\pi t_i) + \epsilon_i$

Prediction

- ▶ Write the equation $x(t_i) = \theta_0 + \theta_1(t_i) + \theta_2(t_i^2) + \theta_3(t_i^3) + \epsilon_i$ in Matrix form

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} 1 & t_1 & t_1^2 & t_1^3 \\ 1 & t_2 & t_2^2 & t_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & t_N & t_N^2 & t_N^3 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

- ▶ Use $\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \boldsymbol{\epsilon}$ and solve for $\boldsymbol{\theta}$

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}'\mathbf{H})^{-1}\mathbf{H}'\mathbf{x}$$

- ▶ Do prediction for any t_i :

$$x(t_i) = \theta_0 + \theta_1(t_i) + \theta_2(t_i^2) + \theta_3(t_i^3) + \epsilon_i$$

Prediction

- ▶ Use any dictionary $\{h_k(t), k = 1, \dots, K\}$
- ▶ Write the equation $x(t_i) = \sum_{k=1}^K \theta_k h_k(t_i) + \epsilon_i$ in matrix form

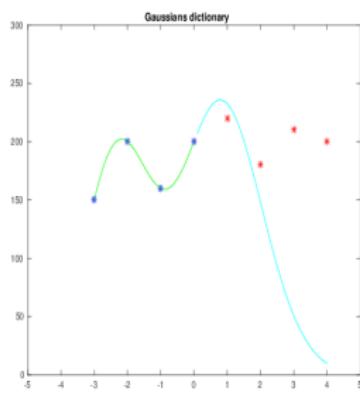
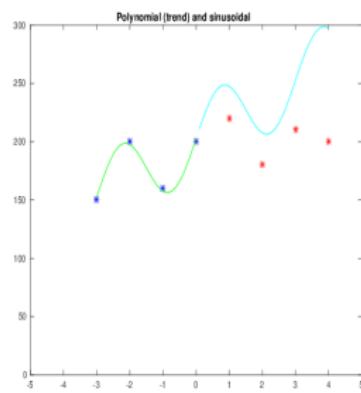
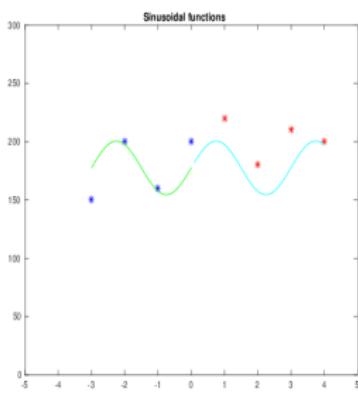
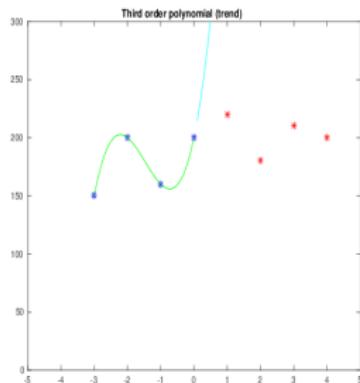
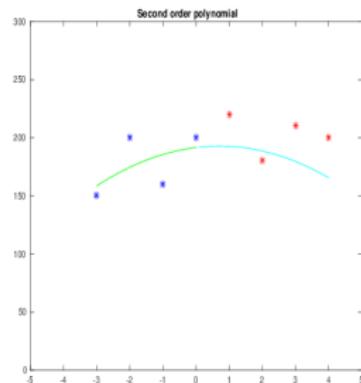
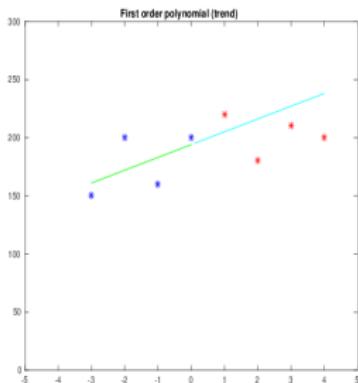
$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} h_1(t_1) & h_2(t_1) & \cdots & h_K(t_1) \\ h_1(t_2) & h_2(t_2) & \cdots & h_K(t_2) \\ h_1(t_3) & h_2(t_3) & \cdots & h_K(t_3) \\ \vdots \\ h_1(t_N) & h_2(t_N) & \cdots & h_K(t_N) \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_K \end{bmatrix}$$

- ▶ Use $\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \boldsymbol{\epsilon}$ and solve for $\boldsymbol{\theta}$

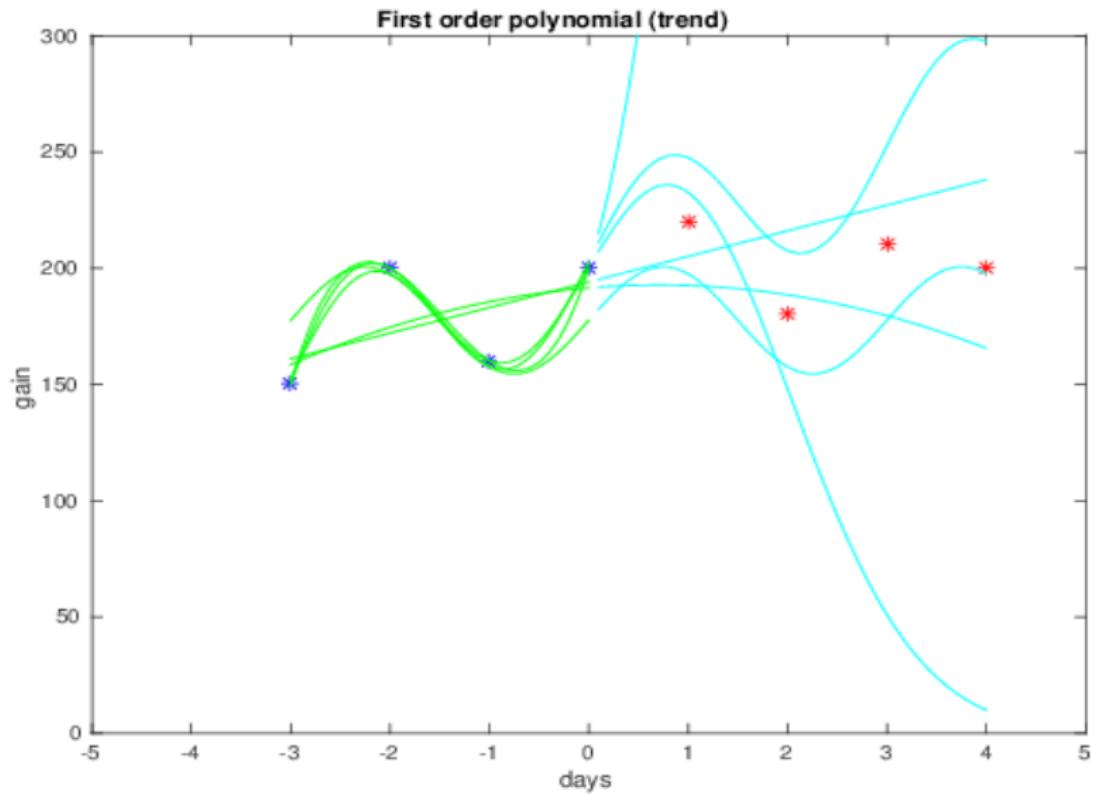
$$\hat{\boldsymbol{\theta}} = (\mathbf{H}'\mathbf{H})^{-1}\mathbf{H}'\mathbf{x}$$

- ▶ Do prediction for any t_i : $x(t_i) = \sum_{k=1}^K \theta_k h_k(t_i) + \epsilon_i$

Prediction examples



Prediction examples

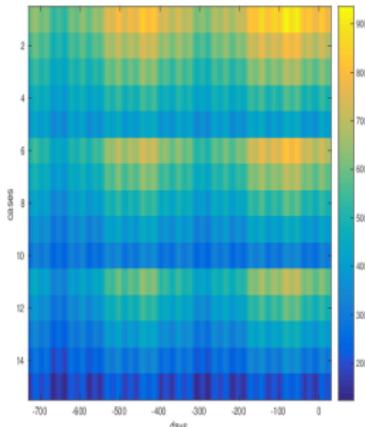
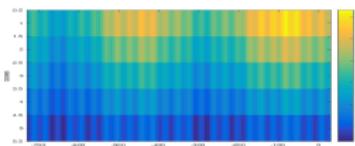
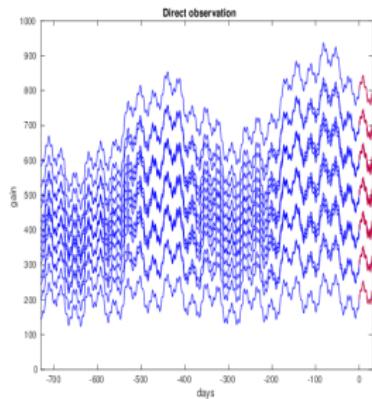
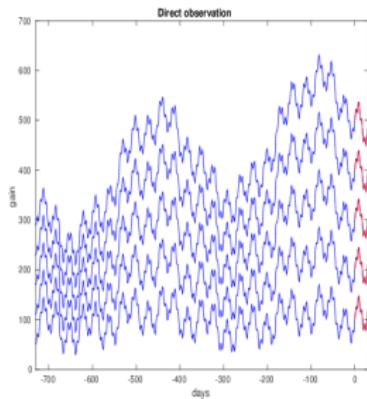
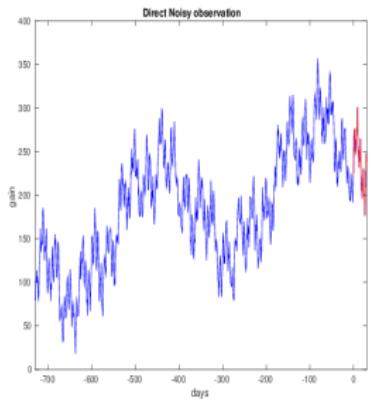


Prediction of gain

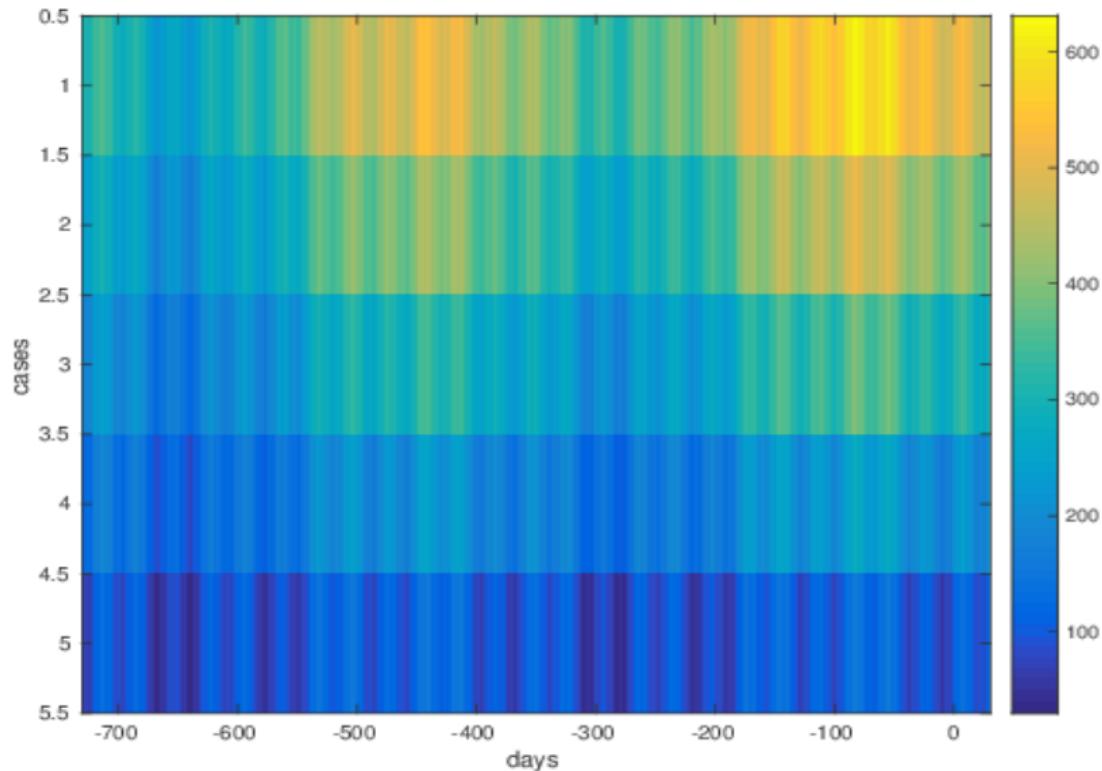
The same problem, but this time you have much more data

- ▶ regular daily since 2 years
- ▶ regular daily but with some missing since 2 years
- ▶ regular daily but with some outliers since 2 years
- ▶ regular daily of yourself and your colleagues with the same rank and positions
- ▶ regular daily of yourself and your colleagues with the same rank and positions, but also other colleagues in your company
- ▶ regular daily of yourself and your colleagues with the same rank and positions, but also other colleagues in your company and in many other companies

Prediction problems examples



Prediction examples



Prediction with direct or indirect observation

The same problem, but this time you want to do daily prediction, but you have the data either weakly or monthly.

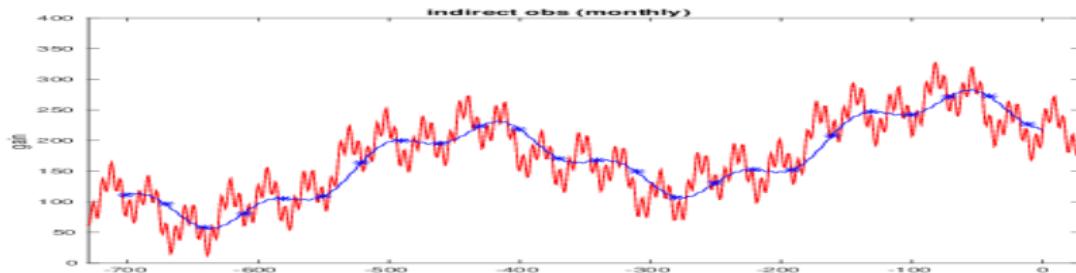
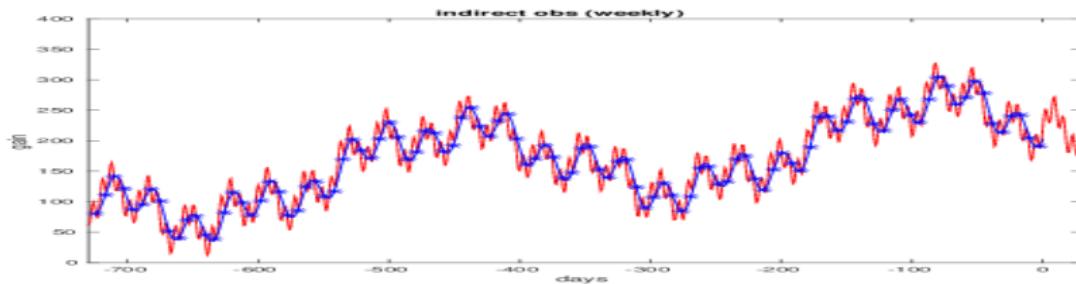
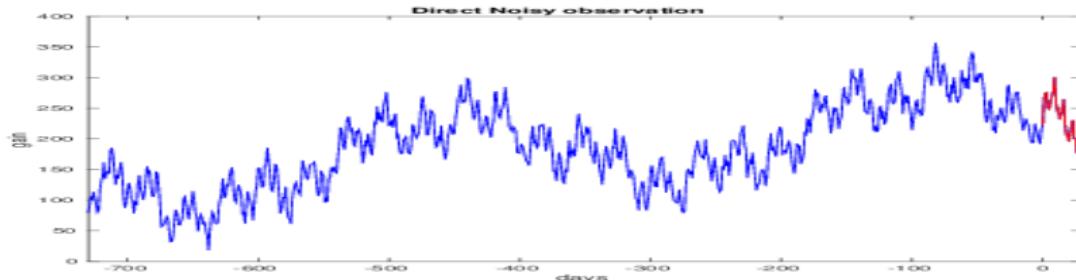
- ▶ Each week, you have the mean value of the week
- ▶ Each month, you have the mean value of the month

If we note $f(n)$ the daily data and $g(m)$ the observed data, then we have

- ▶ last day of the week: $g(m) = f(n = 7 * m)$
- ▶ last day of the month: $g(m) = f(n = 30 * m)$
- ▶ mean value of the week: $g(m) = \sum_{k=1}^7 f(7 * m - k + 1)$
- ▶ mean value of the month: $g(m) = \sum_{k=1}^{30} f(7 * m - k + 1)$
- ▶ Uncertain data

$$g(m) = \sum_{k=1}^K f(7 * m - k + 1) + \epsilon(m)$$

Prediction problems examples



Prediction examples



Population observation, modelling and evolution

We know approximate numbers of population in some of the cities

$$\mathbf{g}(x_i, y_i), i = 1, \dots, M$$

(probably every 4 years since 40 years

$$\mathbf{g}_i(t_n), n = -40 : 4 : 0$$

and we want to know

- ▶ The distribution of the population in the whole country
 $\{\mathbf{f}_{i,j}, i = 1, \dots, I, j = 1, \dots, J\}$
- ▶ The evolution of this distribution year by year
- ▶ Prediction of this distribution of the population in the future years

But also

- ▶ To model the evolution of this distribution and its correlation with some other external events

Population observation, modelling and evolution

Think about modelling:

- Discrete writing

$$\textcolor{blue}{g}(x_i, y_i) = \sum_{(k,l):(x_k, y_l) \in \mathcal{R}_r} \textcolor{red}{f}(x_i - x_k, y_i - y_l) + \epsilon(x_i, y_i)$$

- Continuous writing

$$\textcolor{blue}{g}(x_i, y_i) = \int_{\mathcal{R}_r} \textcolor{red}{f}(x_i - x, y_i - y) dx dy + \epsilon(x_i, y_i)$$

- Differential forms

$$\frac{\partial \textcolor{red}{f}(x, y)}{\partial x} + \textcolor{red}{f}(x, y) = 0 \text{ with initial condition } \textcolor{red}{f}(x, y) = \textcolor{blue}{g}(x, y)$$

Inverse problems scientific communities

Two communities working on Inverse problems:

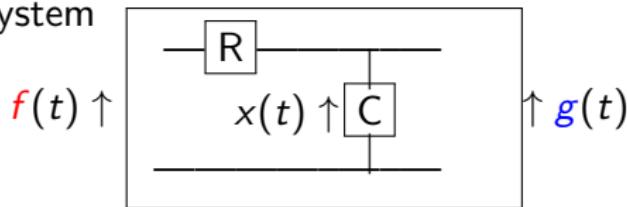
- ▶ Mathematical departments:
Analytical methods: Existence and Uniqueness
Differential equations, PDE
- ▶ Engineering and Computer sciences:
Algebraic methods: Discretization, Uniqueness and Stability
Integral equations, Discretization using Moments method,
Galerkin, ...

Two examples:

- ▶ Deconvolution: Inverse filtering and Wiener filtering
- ▶ X ray Computed Tomography: Radon transform:
Direct Inversion or Filtered Backprojection methods

Differential Equation, State Space and Input-Output

A simple electric system



$$f(t) = R i(t) + v_c(t) = RC \frac{\partial x(t)}{\partial t} + x(t), \quad RC = 1$$

- ▶ Differential Equation Modelling

$$\frac{\partial x(t)}{\partial t} + x(t) = f(t), \quad x(t) = g(t)$$

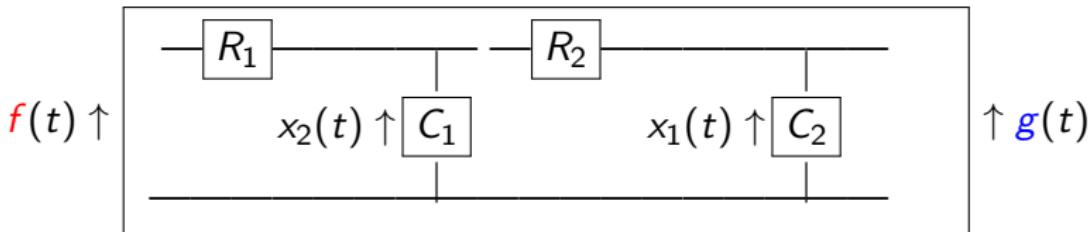
- ▶ State Space Modelling

$$\begin{cases} \frac{\partial x(t)}{\partial t} = -x(t) + f(t) \\ g(t) = x(t) \end{cases}$$

- ▶ Input-Output Modelling

$$\begin{cases} \frac{\partial x(t)}{\partial t} = -x(t) + f(t) \\ g(t) = x(t) \end{cases} \rightarrow \begin{cases} pX(p) = -X(p) + F(p) \rightarrow X(p) = \frac{1}{p+1}F(p) \\ g(t) = x(t) = h(t) * f(t), \quad h(t) = \exp[-t] \end{cases}$$

A more complex electric system example



$$f(t) = \frac{\partial x_2(t)}{\partial t} + x_2(t), \quad x_2(t) = \frac{\partial x_1(t)}{\partial t} + x_1(t)$$

- ▶ Differential Equation model: $\frac{\partial^2 x_1(t)}{\partial t^2} + 2\frac{\partial x_1(t)}{\partial t} + x_1(t) = f(t)$
- ▶ State space model

$$\begin{cases} \begin{bmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} f(t) \\ g(t) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} x_1(t) \end{cases}$$

- ▶ Input-Output Model: $g(t) = h(t) * f(t)$

Design/Control Inverse problems examples

Simple Electrical system:

$$a \frac{\partial x(t)}{\partial t} + x(t) = f(t), \quad x(0) = x_0, \quad g(t) = x(t)$$

- ▶ Design: $\theta = a = RC$
 - ▶ Forward: Given $\theta = a$ and $f(t), t > 0$, find $x(t), t > 0$
 - ▶ Inverse: Given $x(t)$ and $f(t)$ find $\theta = a$
- ▶ Control: $f(t)$
 - ▶ Forward: Given $\theta = a$ and $f(t), t > 0$, find $x(t), t > 0$
 - ▶ Inverse: Given $\theta = a$ and $x(t), t > 0$, find $f(t)$

More complex Electrical system:

$$f(t) = b \frac{\partial x_2(t)}{\partial t} + x_2(t), \quad x_2(t) = a \frac{\partial x_1(t)}{\partial t} + x_1(t), \quad g(t) = x_1(t)$$

$$\theta = (a = R_1 C_1, b = R_2 C_2)$$

Design/Control Inverse problems examples

Mass-spring-dashpot system

$$m \frac{\partial^2 x(t)}{\partial t^2} + c \frac{\partial x(t)}{\partial t} + k = F(t), \quad x(0) = x_0, \quad \frac{\partial x}{\partial t}(0) = v_0$$

- ▶ Design: $\theta = (m, c, k)$
 - ▶ Forward: Given $\theta = (m, c, k)$, x_0, v_0 and $F(t), t > 0$, find $x(t), t > 0$
 - ▶ Inverse: Given $x(t)$ for $t > 0$, v_0 , $F(t)$ find $\theta = (m, c, k)$
- ▶ Control: $F(t)$
 - ▶ Forward: Given $\theta = (m, c, k)$, x_0, v_0 and $F(t), t > 0$, find $x(t), t > 0$
 - ▶ Inverse: Given $\theta = (m, c, k)$, v_0 and $x(t), t > 0$, find $F(t)$

Input-Output model

- ▶ Linear Systems

- ▶ Single Input Single Output (SISO) systems

$$y(t) = \int h(t, \tau) u(\tau) d\tau$$

- ▶ Multi Input Multi Output (MIMO) systems

$$\mathbf{y}(t) = \int \mathbf{H}(t, \tau) \mathbf{u}(\tau) d\tau$$

- ▶ Linear Time Invariant System

- ▶ SISO Convolution

$$y(t) = h(t) * u(t) = \int h(t - \tau) u(\tau) d\tau$$

- ▶ MIMO Convolution

$$\mathbf{y}(t) = \int \mathbf{H}(t - \tau) \mathbf{u}(\tau) d\tau$$

- ▶ Impulse response $h(t)$ or $\mathbf{H}(t) = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & h_{ij}(t) & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$

State space model: Continuous case

Dynamic systems:

- ▶ Single Input Single Output (SISO) system:

$$\begin{cases} \dot{\mathbf{x}}(t) = A\mathbf{x}(t) + B\mathbf{u}(t) & \text{State equation} \\ \mathbf{y}(t) = C\mathbf{x}(t) + D\mathbf{v}(t) & \text{Observation equation} \end{cases}$$

- ▶ Multiple Input Multiple Output (MIMO) system:

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{H}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t) & \text{State equation} \\ \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{v}(t) & \text{Observation equation} \end{cases}$$

H, **B**, **C** and **D** are the matrices of the system.

Modelling with Partial Differential Equations

- ▶ Different PDE

$$\frac{\partial^2 \mathbf{f}(x, y)}{\partial x^2} + \frac{\partial^2 \mathbf{f}(x, y)}{\partial y^2} + \mathbf{f}(x, y) = 0$$

$$\frac{\partial^2 \mathbf{f}(x, y)}{\partial x^2} + \frac{\partial^2 \mathbf{f}(x, y)}{\partial y^2} + \frac{\partial \mathbf{f}(x, y)}{\partial x} \frac{\partial \mathbf{f}(x, y)}{\partial y} \mathbf{f}(x, y) = 0$$

with initial cond. $\mathbf{f}(x, y) = \mathbf{g}(x, y)$

Prediction with indirect observation

- ▶ Data available every K days with uncertainty

$$\textcolor{blue}{g}(m) = \sum_{k=1}^K \textcolor{red}{f}(K * m - k + 1) + \epsilon(m)$$

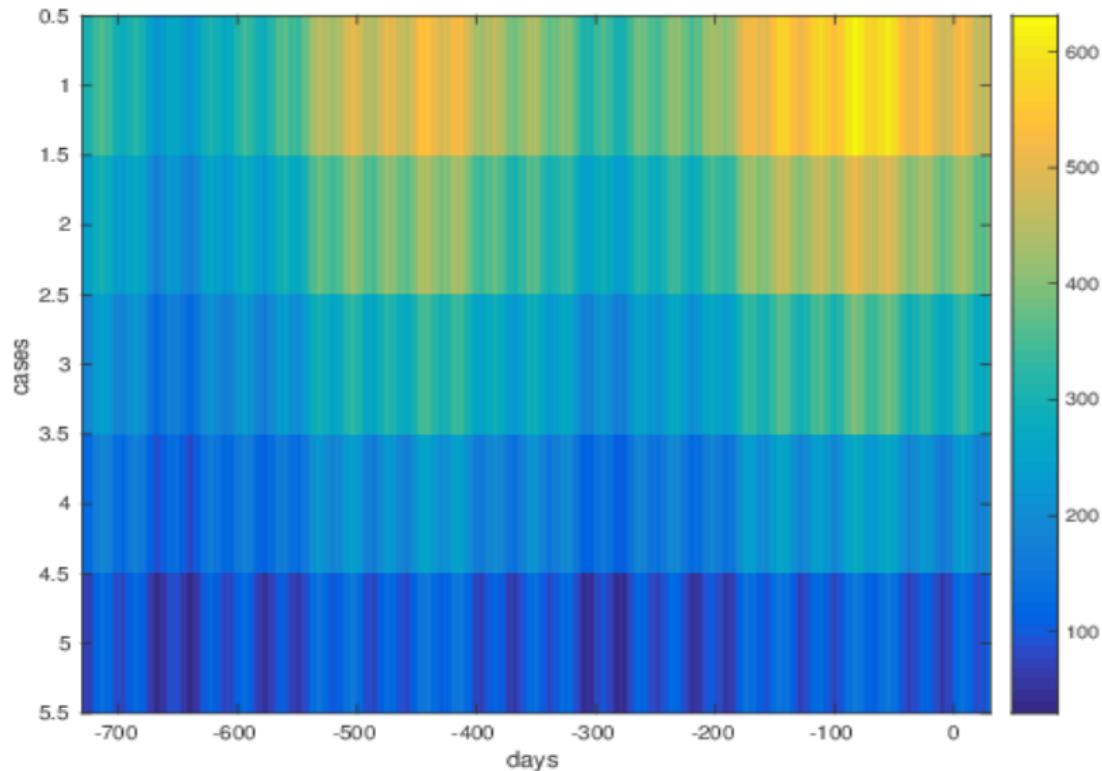
- ▶ more general: Convolution

$$\textcolor{blue}{g}(m) = \sum_{k=1}^K h(k) \textcolor{red}{f}(n - k + 1) + \epsilon(m)$$

- ▶ One can show easily that both can be written as

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}$$

Prediction examples



Simple Examples

f_1	f_5	f_9	f_{13}	g_8	f_{11}	f_{12}	f_{13}	f_{14}	g_{24}
f_2	f_6	f_{10}	f_{14}	g_7	f_{21}	f_{22}	f_{23}	f_{24}	g_{23}
f_3	f_7	f_{11}	f_{15}	g_6	f_{31}	f_{32}	f_{33}	f_{34}	g_{22}
f_4	f_8	f_{12}	f_{16}	g_5	f_{41}	f_{42}	f_{43}	f_{44}	g_{21}
g_1	g_2	g_3	g_4		g_{11}	g_{12}	g_{13}	g_{14}	

Noting also by:

$$\mathbf{g}_1 = [g_1, \dots, g_4]^t = [g_{11}, \dots, g_{14}]^t,$$

$$\mathbf{g}_2 = [g_5, \dots, g_8]^t = [g_{21}, \dots, g_{24}]^t$$

and the matrices \mathbf{H}_1 , \mathbf{H}_2 and \mathbf{H} such that:

$$\mathbf{g}_1 = \mathbf{H}_1 \mathbf{f}, \quad \mathbf{g}_2 = \mathbf{H}_2 \mathbf{f}, \quad \mathbf{g} = \mathbf{H} \mathbf{f} = \begin{bmatrix} \mathbf{H}_1 \\ \mathbf{H}_2 \end{bmatrix} \mathbf{f}$$

Inversion, Generalized inversion

- ▶ Forward problem: Given \mathbf{f} compute \mathbf{g} :

$$f = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \begin{cases} \mathbf{g}_1 = \mathbf{H}_1\mathbf{f} = [0 \ 2 \ 2 \ 0]^t, \\ \mathbf{g}_2 = \mathbf{H}_2\mathbf{f} = [0 \ 2 \ 2 \ 0]^t, \\ \mathbf{g} = \begin{bmatrix} \mathbf{g}_1 \\ \mathbf{g}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{H}_1 \\ \mathbf{H}_2 \end{bmatrix} \mathbf{f} \\ \mathbf{g} = \mathbf{H}\mathbf{f} = [0 \ 2 \ 2 \ 0 \ 0 \ 2 \ 2 \ 0]^t \end{cases}$$

- ▶ Inverse problem: Given \mathbf{g} find \mathbf{f} .

Many possible solutions:

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} -.5 & 0 & 0 & .5 \\ 1 & 2 & 0 & -1 \\ -1 & 0 & 2 & 1 \\ 0.5 & 0 & 0 & -.5 \end{bmatrix} \quad \begin{bmatrix} -.5 & 0 & 0 & .5 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ .5 & 0 & 0 & -.5 \end{bmatrix}$$

MN, LS and MNLS solutions

- ▶ Minimum Norme (MN) solution:

$$\widehat{\mathbf{f}} = \underset{\mathbf{H}\mathbf{f}=\mathbf{g}}{\arg \min} \left\{ \|\mathbf{f}\|^2 \right\}$$

and if $\mathbf{H}\mathbf{H}^t$ was invertible, then we had: $\widehat{\mathbf{f}} = \mathbf{H}^t(\mathbf{H}\mathbf{H}^t)^{-1}\mathbf{g}$.
But

$$\text{svd}(\mathbf{H}\mathbf{H}^t) = [8 \ 4 \ 4 \ 4 \ 4 \ 4 \ 4 \ 0]$$

- ▶ Least Squares (LS) solution:

$$\widehat{\mathbf{f}} = \underset{\mathbf{f}}{\arg \min} \left\{ \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 \right\},$$

and if $\mathbf{H}^t\mathbf{H}$ was invertible, then we had: $\widehat{\mathbf{f}} = (\mathbf{H}^t\mathbf{H})^{-1}\mathbf{H}^t\mathbf{g}$.
But

$$\text{svd}(\mathbf{H}^t\mathbf{H}) = [8 \ 4 \ 4 \ 4 \ 4 \ 4 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$$

SVD and MNLS solutions

- ▶ Truncation of the singular values to define an unique Generalized Inverse solution

$$\hat{\mathbf{f}} = \sum_{k=1}^k \frac{<\mathbf{g}, \mathbf{u}_k>}{\lambda_k} \mathbf{v}_k$$

where \mathbf{u}_k and \mathbf{v}_k are, respectively, the eigenvectors of $\mathbf{A}\mathbf{A}^t$ and $\mathbf{A}^t\mathbf{A}$ and λ_k their corresponding eigen values.

$$\hat{\mathbf{f}} = \begin{bmatrix} -0.2500 & 0.2500 & 0.2500 & -0.2500 \\ 0.2500 & 0.7500 & 0.7500 & 0.2500 \\ 0.2500 & 0.7500 & 0.7500 & 0.2500 \\ -0.2500 & 0.2500 & 0.2500 & -0.2500 \end{bmatrix}$$

- ▶ MNLS

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \left\{ \|\mathbf{g} - \mathbf{Hf}\|^2 + \|\mathbf{f}\|^2 \right\},$$

Regularization theory

Inverse problems = III posed problems

→ Need for prior information

Functional space (Tikhonov):

$$\mathbf{g} = \mathcal{H}(\mathbf{f}) + \epsilon \rightarrow J(\mathbf{f}) = \|\mathbf{g} - \mathcal{H}(\mathbf{f})\|_2^2 + \lambda \|\mathcal{D}\mathbf{f}\|_2^2$$

Finite dimensional space (Philips & Towney): $\mathbf{g} = \mathbf{H}(\mathbf{f}) + \epsilon$

- Minimum norme LS (MNLS): $J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}(\mathbf{f})\|^2 + \lambda \|\mathbf{f}\|^2$
- Classical regularization: $J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}(\mathbf{f})\|^2 + \lambda \|\mathbf{Df}\|^2$
- More general regularization:

$$J(\mathbf{f}) = Q(\mathbf{g} - \mathbf{H}(\mathbf{f})) + \lambda \Omega(\mathbf{Df})$$

or

$$J(\mathbf{f}) = \Delta_1(\mathbf{g}, \mathbf{H}(\mathbf{f})) + \lambda \Delta_2(\mathbf{f}, \mathbf{f}_\infty)$$

Limitations:

- Errors are considered implicitly white and Gaussian
- Limited prior information on the solution
- Lack of tools for the determination of the hyperparameters

Basic Bayes

- ▶ Product rules

$$P(A, B) = P(A|B)P(B) = P(B|A)P(A) \rightarrow P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- ▶ Sum rule

$$P(B) = \sum_A P(A|B)P(B)$$

- ▶ Bayes rule (discrete events)

$$P(A|B) = \frac{P(B|A)P(A)}{\sum_A P(B|A)P(A)}$$

- ▶ $P(\text{hypothesis}|\text{data}) = \frac{P(\text{data}|\text{hypothesis})P(\text{hypothesis})}{P(\text{data})}$
- ▶ Bayes rule (Continuous variables with finite parametric models)

$$p(\boldsymbol{\theta}|\mathbf{d}) = \frac{p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{d})} = \frac{p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}} \propto p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta})$$

Basic Bayes

- ▶ $P(\text{hypothesis}|\text{data}) = \frac{P(\text{data}|\text{hypothesis})P(\text{hypothesis})}{P(\text{data})}$
- ▶ Bayes rule tells us how to do inference about hypotheses from data.
- ▶ Finite parametric models:

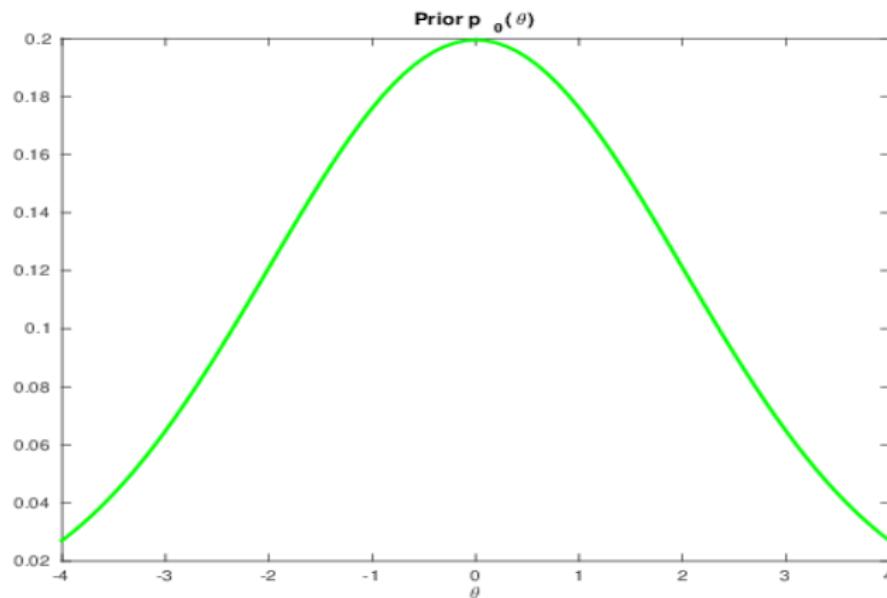
$$p(\boldsymbol{\theta}|\mathbf{d}) = \frac{p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{d})}$$

- ▶ Forward model: $p(\mathbf{d}|\boldsymbol{\theta})$
called also likelihood of the parameters in data $\mathcal{L}(\boldsymbol{\theta}) = p(\mathbf{d}|\boldsymbol{\theta})$
- ▶ Prior knowledge: $p(\boldsymbol{\theta})$
- ▶ Posterior knowledge: $p(\boldsymbol{\theta}|\mathbf{d})$
- ▶ Evidence: $p(\mathbf{d}) = \int p(\mathbf{d}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$

Bayesian inference: simple one parameter case

$$p(\theta), \mathcal{L}(\theta) = p(\mathbf{d}|\theta) \longrightarrow p(\theta|\mathbf{d}) \propto \mathcal{L}(\theta) p(\theta)$$

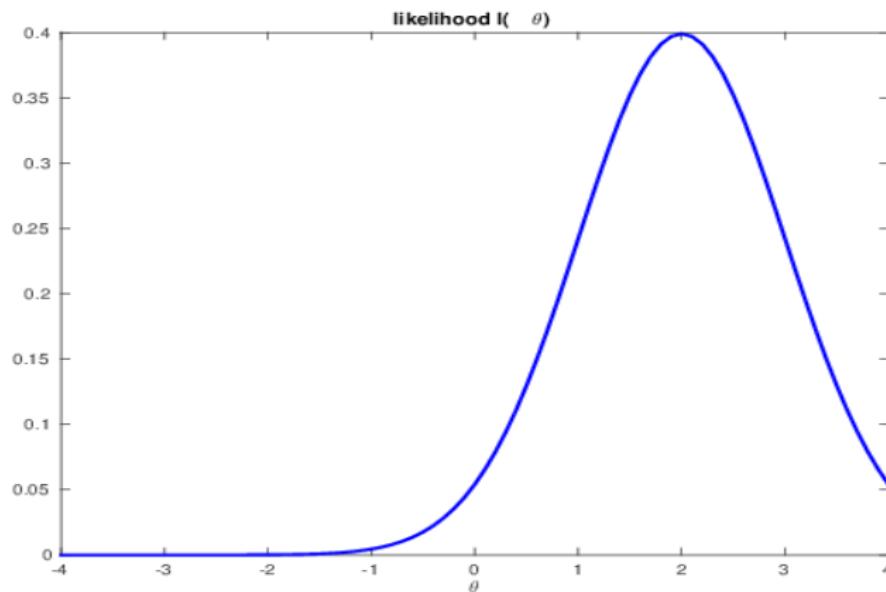
Prior: $p(\theta)$



Bayesian inference: simple one parameter case

$$p(\theta), \mathcal{L}(\theta) = p(\mathbf{d}|\theta) \longrightarrow p(\theta|\mathbf{d}) \propto \mathcal{L}(\theta) p(\theta)$$

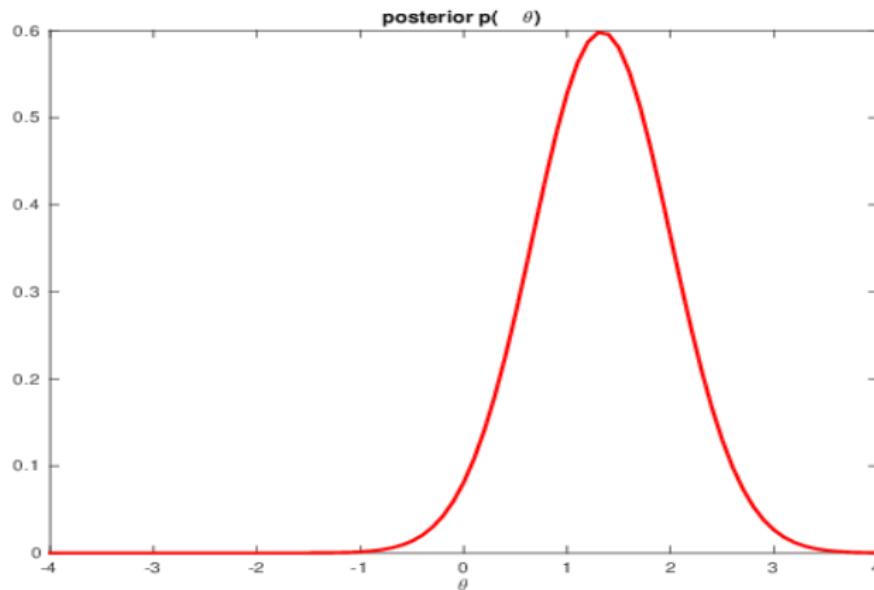
Likelihood: $\mathcal{L}(\theta) = p(\mathbf{d}|\theta)$



Bayesian inference: simple one parameter case

$$p(\theta), \mathcal{L}(\theta) = p(\mathbf{d}|\theta) \longrightarrow p(\theta|\mathbf{d}) \propto \mathcal{L}(\theta) p(\theta)$$

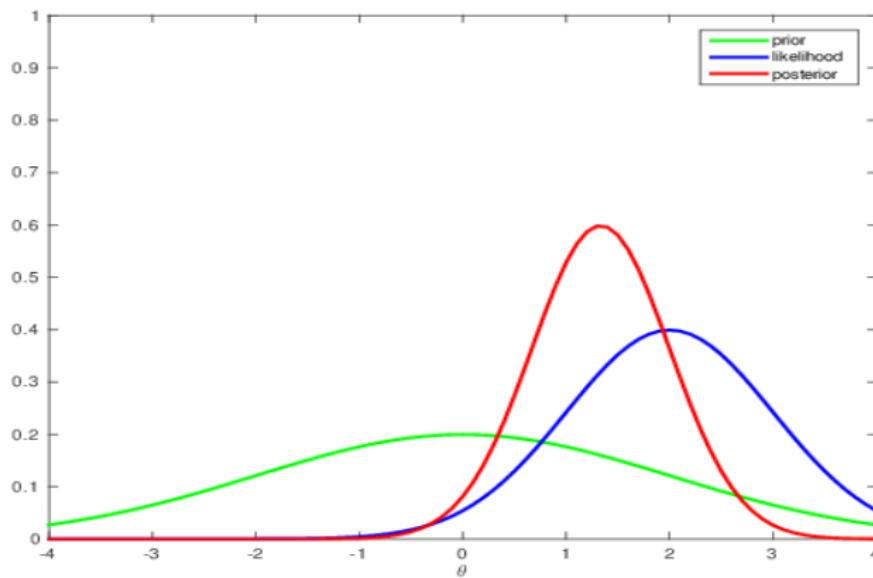
Posterior: $p(\theta|\mathbf{d}) \propto p(\mathbf{d}|\theta) p(\theta)$



Bayesian inference: simple one parameter case

$$p(\theta), \mathcal{L}(\theta) = p(\mathbf{d}|\theta) \longrightarrow p(\theta|\mathbf{d}) \propto \mathcal{L}(\theta) p(\theta)$$

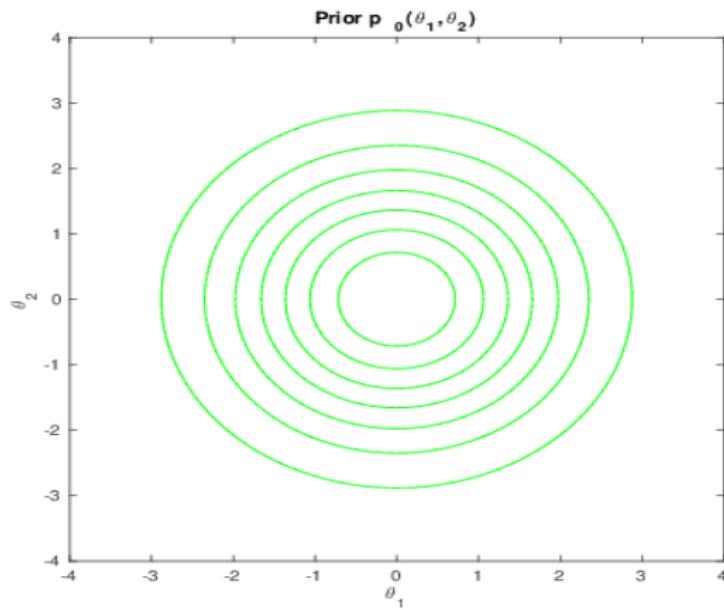
Prior, Likelihood and Posterior:



Bayesian inference: simple two parameters case

$$p(\theta_1, \theta_2), \mathcal{L}(\theta_1, \theta_2) = p(\mathbf{d}|\theta_1, \theta_2) \longrightarrow p(\theta_1, \theta_2|\mathbf{d}) \propto \mathcal{L}(\theta_1, \theta_2) p(\theta_1, \theta_2)$$

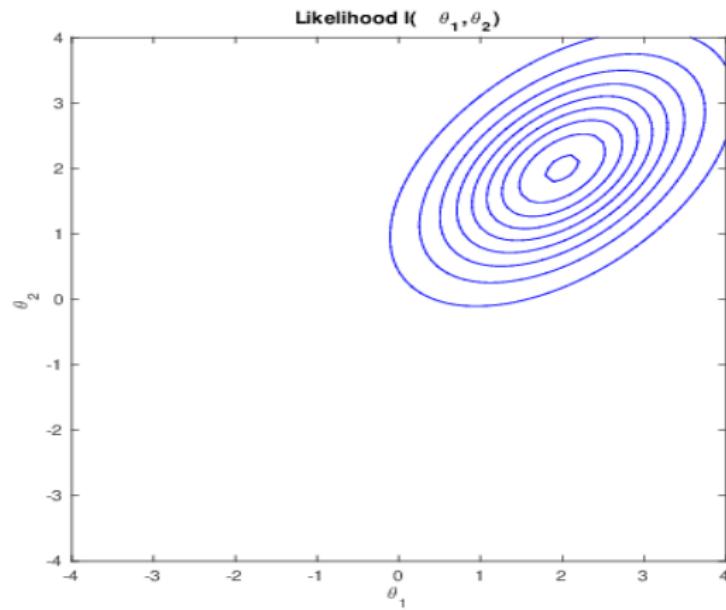
Prior: $p(\theta_1, \theta_2)$



Bayesian inference: simple one parameter case

$$p(\theta_1, \theta_2), \mathcal{L}(\theta_1, \theta_2) = p(\mathbf{d}|\theta_1, \theta_2) \longrightarrow p(\theta_1, \theta_2|\mathbf{d}) \propto \mathcal{L}(\theta_1, \theta_2) p(\theta_1, \theta_2)$$

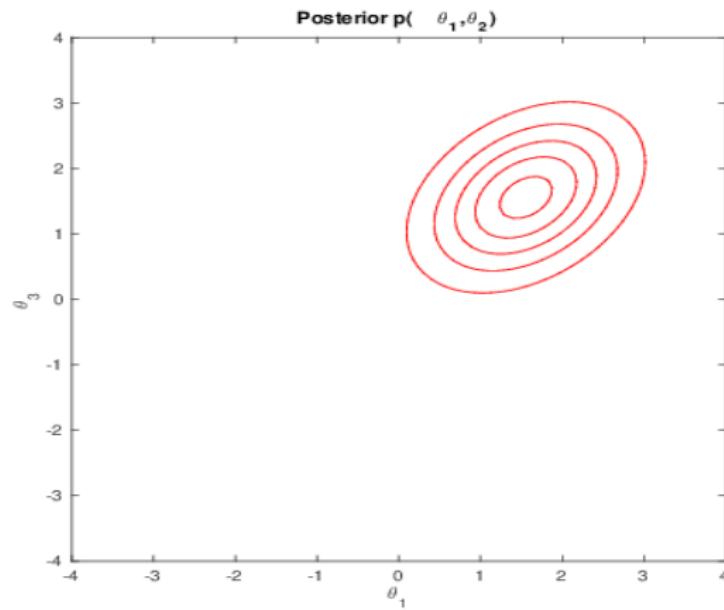
Likelihood: $\mathcal{L}(\theta_1, \theta_2) = p(\mathbf{d}|\theta_1, \theta_2)$



Bayesian inference: simple one parameter case

$$p(\theta_1, \theta_2), \mathcal{L}(\theta_1, \theta_2) = p(\mathbf{d}|\theta_1, \theta_2) \longrightarrow p(\theta_1, \theta_2|\mathbf{d}) \propto \mathcal{L}(\theta_1, \theta_2) p(\theta_1, \theta_2)$$

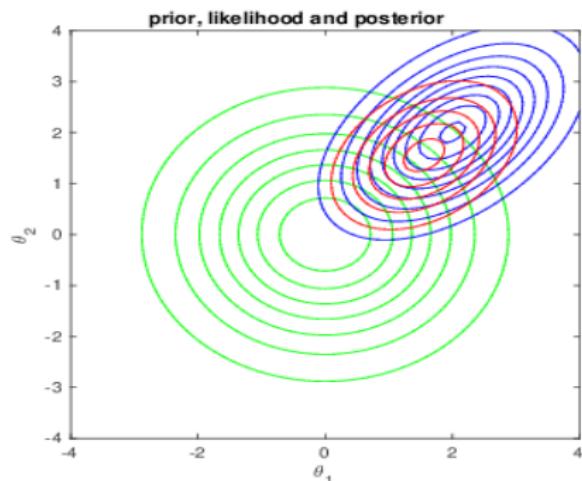
Posterior: $p(\theta_1, \theta_2|\mathbf{d}) \propto p(\mathbf{d}|\theta_1, \theta_2) p(\theta_1, \theta_2)$



Bayesian inference: simple one parameter case

$$p(\theta_1, \theta_2), \mathcal{L}(\theta_1, \theta_2) = p(\mathbf{d}|\theta_1, \theta_2) \longrightarrow p(\theta_1, \theta_2|\mathbf{d}) \propto \mathcal{L}(\theta_1, \theta_2) p(\theta_1, \theta_2)$$

Prior, Likelihood and Posterior:



Bayes: 1D case

$$p(\theta|\mathbf{d}) = \frac{p(\mathbf{d}|\theta) p(\theta)}{p(\mathbf{d})} \propto p(\mathbf{d}|\theta) p(\theta)$$

- ▶ Maximum A Posteriori (MAP)

$$\hat{\theta} = \arg \max_{\theta} \{p(\theta|\mathbf{d})\} = \arg \max_{\theta} \{p(\mathbf{d}|\theta) p(\theta)\}$$

- ▶ Posterior Mean

$$\hat{\theta} = E_{p(\theta|\mathbf{d})} \{\theta\} = \int \theta p(\theta|\mathbf{d}) d\theta$$

- ▶ Region of high probabilities

$$[\hat{\theta}_1, \hat{\theta}_2] : \int_{\hat{\theta}_1}^{\hat{\theta}_2} p(\theta|\mathbf{d}) d\theta = 1 - \alpha$$

- ▶ Sampling and exploring

$$\theta \sim p(\theta|\mathbf{d})$$

Bayesian inference: great dimensional case

- ▶ Simple Linear case: $\mathbf{d} = \mathbf{H}\theta + \epsilon$
- ▶ Gaussian priors:

$$p(\mathbf{d}|\theta) = \mathcal{N}(\mathbf{d}|\mathbf{H}\theta, v_\epsilon \mathbf{I})$$
$$p(\theta) = \mathcal{N}(\theta|0, v_\theta \mathbf{I})$$

- ▶ Gaussian posterior:

$$p(\theta|\mathbf{d}) = \mathcal{N}(\theta|\hat{\theta}, \hat{\mathbf{V}})$$
$$\hat{\theta} = [\mathbf{H}'\mathbf{H} + \lambda \mathbf{I}]^{-1} \mathbf{H}'\mathbf{d}, \quad \lambda = \frac{v_\epsilon}{v_\theta}$$
$$\hat{\mathbf{V}} = [\mathbf{H}'\mathbf{H} + \lambda \mathbf{I}]^{-1}$$

- ▶ Computation of $\hat{\theta}$ can be done via optimization of:
$$J(\theta) = -\ln p(\theta|\mathbf{d}) = \frac{1}{2v_\epsilon} \|\mathbf{d} - \mathbf{H}\theta\|^2 + \frac{1}{2v_\theta} \|\theta\|^2 + c$$
- ▶ Computation of $\hat{\mathbf{V}} = [\mathbf{H}'\mathbf{H} + \lambda \mathbf{I}]^{-1}$ needs great dimensional matrix inversion.

Bayesian inference: great dimensional case

- ▶ Gaussian posterior:

$$p(\boldsymbol{\theta}|\mathbf{d}) = \mathcal{N}(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}, \hat{\mathbf{V}}),$$
$$\hat{\boldsymbol{\theta}} = [\mathbf{H}'\mathbf{H} + \lambda\mathbf{I}]^{-1}\mathbf{H}'\mathbf{d}, \quad \hat{\mathbf{V}} = [\mathbf{H}'\mathbf{H} + \lambda\mathbf{I}]^{-1}, \quad \lambda = \frac{v_\epsilon}{v_\theta}$$

- ▶ Computation of $\hat{\boldsymbol{\theta}}$ can be done via optimization of:

$$J(\boldsymbol{\theta}) = -\ln p(\boldsymbol{\theta}|\mathbf{d}) = c + \|\mathbf{d} - \mathbf{H}\boldsymbol{\theta}\|^2 + \lambda\|\boldsymbol{\theta}\|^2$$

- ▶ Gradient based methods:

$$\nabla J(\boldsymbol{\theta}) = -2\mathbf{H}'(\mathbf{d} - \mathbf{H}\boldsymbol{\theta}) + 2\lambda\boldsymbol{\theta}$$

- ▶ constant step, Steepest descend, ...

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \alpha^{(k)} \nabla J(\boldsymbol{\theta}^{(k)}) = \boldsymbol{\theta}^{(k)} + 2\alpha^{(k)} [\mathbf{H}'(\mathbf{d} - \mathbf{H}\boldsymbol{\theta}) + \lambda\boldsymbol{\theta}]$$

- ▶ Conjugate Gradient, ...

- ▶ At each iteration, we need to be able to compute:

- ▶ Forward operation: $\hat{\mathbf{d}} = \mathbf{H}\boldsymbol{\theta}^{(k)}$

- ▶ Backward (Adjoint) operation: $\mathbf{H}^t(\mathbf{d} - \hat{\mathbf{d}})$

Bayesian inference: great dimensional case

- ▶ Computation of $\hat{\mathbf{V}} = [\mathbf{H}'\mathbf{H} + \lambda\mathbf{I}]^{-1}$ needs great dimensional matrix inversion.
- ▶ Almost impossible except in particular cases of Toeplitz, Circulante, TBT, CBC,... where we can diagonalize it via Fast Fourier Transform (FFT).
- ▶ Recursive use of the data and recursive update of $\hat{\boldsymbol{\theta}}$ and $\hat{\mathbf{V}}$ leads to Kalman Filtering which are still computationally demanding for High dimensional data.
- ▶ We also need to generate samples from this posterior: There are many special sampling tools.
- ▶ Mainly two categories: Using the covariance matrix \mathbf{V} or its inverse (Precision matrix) $\Lambda = \mathbf{V}^{-1}$

Bayesian inference: non Gaussian priors case

- ▶ Linear forward model: $\mathbf{d} = \mathbf{H}\boldsymbol{\theta} + \boldsymbol{\epsilon}$
- ▶ Gaussian noise model:

$$p(\mathbf{d}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{d}|\mathbf{H}\boldsymbol{\theta}, v_\epsilon \mathbf{I}) \propto \exp\left[-\frac{1}{2v_\epsilon} \|\mathbf{d} - \mathbf{H}\boldsymbol{\theta}\|_2^2\right]$$

- ▶ Sparsity enforcing prior:

$$p(\boldsymbol{\theta}) \propto \exp[\alpha \|\boldsymbol{\theta}\|_1]$$

- ▶ Posterior:

$$p(\boldsymbol{\theta}|\mathbf{d}) \propto \exp\left[-\frac{1}{2v_\epsilon} J(\boldsymbol{\theta})\right] \text{ with } J(\boldsymbol{\theta}) = \|\mathbf{d} - \mathbf{H}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1, \lambda = 2v_\epsilon\alpha$$

- ▶ Computation of $\hat{\boldsymbol{\theta}}$ can be done via optimization of $J(\boldsymbol{\theta})$
- ▶ Other computations are much more difficult.

Bayes Rule for Machine Learning (Simple case)

- Inference on the parameters: Learning from data \mathbf{d} :

$$p(\boldsymbol{\theta}|\mathbf{d}, \mathcal{M}) = \frac{p(\mathbf{d}|\boldsymbol{\theta}, \mathcal{M}) p(\boldsymbol{\theta}|\mathcal{M})}{p(\mathbf{d}|\mathcal{M})}$$

- Model Comparison:

$$p(\mathcal{M}_k|\mathbf{d}) = \frac{p(\mathbf{d}|\mathcal{M}_k) p(\mathcal{M}_k)}{p(\mathbf{d})}$$

with

$$p(\mathbf{d}|\mathcal{M}_k) = \int p(\mathbf{d}|\boldsymbol{\theta}, \mathcal{M}_k) p(\boldsymbol{\theta}|\mathcal{M}) d\boldsymbol{\theta}$$

- Prediction with selected model:

$$p(\mathbf{z}|\mathcal{M}_k) = \int p(\mathbf{z}|\boldsymbol{\theta}, \mathcal{M}_k) p(\boldsymbol{\theta}|\mathbf{d}, \mathcal{M}_k) d\boldsymbol{\theta}$$

Approximation methods

- ▶ Laplace approximation
- ▶ Bayesian Information Criterion (BIC)
- ▶ Variational Bayesian Approximations (VBA)
- ▶ Expectation Propagation (EP)
- ▶ Markov chain Monte Carlo methods (MCMC)
- ▶ Exact Sampling

Laplace Approximation

- ▶ Data set \mathbf{d} , models $\mathcal{M}_1, \dots, \mathcal{M}_K$, parameters $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K$
- ▶ Model Comparison:

$$p(\boldsymbol{\theta}, \mathbf{d} | \mathcal{M}) = p(\mathbf{d} | \boldsymbol{\theta}, \mathcal{M}) p(\boldsymbol{\theta} | \mathcal{M})$$

$$p(\boldsymbol{\theta} | \mathbf{d}, \mathcal{M}) = p(\boldsymbol{\theta}, \mathbf{d} | \mathcal{M}) / p(\mathbf{d} | \mathcal{M})$$

$$p(\mathbf{d} | \mathcal{M}) = \int p(\mathbf{d} | \boldsymbol{\theta}, \mathcal{M}) p(\boldsymbol{\theta} | \mathcal{M}) d\boldsymbol{\theta}$$

- ▶ For large amount of data (relative to number of parameters, m), $p(\boldsymbol{\theta} | \mathbf{d}, \mathcal{M})$ is approximated by a Gaussian around its maximum (MAP) $\hat{\boldsymbol{\theta}}$:

$$p(\boldsymbol{\theta} | \mathbf{d}, \mathcal{M}) \approx (2\pi)^{-m/2} |\mathbf{A}|^{1/2} \exp \left[-\frac{1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})' \mathbf{A} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right]$$

where $A_{i,j} = \frac{d^2}{\theta_i \theta_j} \ln p(\boldsymbol{\theta} | \mathbf{d}, \mathcal{M})$ is the $m \times m$ Hessian matrix.

- ▶ $p(\mathbf{d} | \mathcal{M}) = p(\boldsymbol{\theta}, \mathbf{d} | \mathcal{M}) / p(\boldsymbol{\theta} | \mathbf{d}, \mathcal{M})$ and evaluating it at $\hat{\boldsymbol{\theta}}$:

$$\ln p(\mathbf{d} | \mathcal{M}_k) \approx \ln p(\mathbf{d} | \hat{\boldsymbol{\theta}}, \mathcal{M}_k) + \ln p(\hat{\boldsymbol{\theta}} | \mathcal{M}_k) + \frac{m}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|$$

- ▶ Needs computation of $\hat{\boldsymbol{\theta}}$ and $|\mathbf{A}|$.

Bayesian Information Criteria (BIC)

- ▶ BIC is obtained from the Laplace approximation

$$\ln p(\mathbf{d}|\mathcal{M}_k) \approx \ln p(\hat{\boldsymbol{\theta}}|\mathcal{M}_k) + p(\mathbf{d}|\hat{\boldsymbol{\theta}}, \mathcal{M}_k) + \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\mathbf{A}|$$

by taking the large sample limit ($n \mapsto \infty$) where n is the number of data points:

$$\ln p(\mathbf{d}|\mathcal{M}_k) \approx p(\mathbf{d}|\hat{\boldsymbol{\theta}}, \mathcal{M}_k) - \frac{d}{2} \ln(n)$$

- ▶ Easy to compute
- ▶ It does not depend on the prior
- ▶ It is equivalent to MDL criterion
- ▶ Assumes that as ($n \mapsto \infty$), all the parameters are identifiable.
- ▶ Danger: counting parameters can be deceiving (sinusoid, infinite dim models)

Bayes Rule for Machine Learning with hidden variables

- ▶ Data: \mathbf{d} , Hidden Variables: \mathbf{x} , Parameters: $\boldsymbol{\theta}$, Model: \mathcal{M}
- ▶ Bayes rule

$$p(\mathbf{x}, \boldsymbol{\theta} | \mathbf{d}, \mathcal{M}) = \frac{p(\mathbf{d} | \mathbf{x}, \boldsymbol{\theta}, \mathcal{M}) p(\mathbf{x} | \boldsymbol{\theta}, \mathcal{M}) p(\boldsymbol{\theta} | \mathcal{M})}{p(\mathbf{d} | \mathcal{M})}$$

- ▶ Model Comparison

$$p(\mathcal{M}_k | \mathbf{d}) = \frac{p(\mathbf{d} | \mathcal{M}_k) p(\mathcal{M}_k)}{p(\mathbf{d})}$$

with

$$p(\mathbf{d} | \mathcal{M}_k) = \int \int p(\mathbf{d} | \mathbf{x}, \boldsymbol{\theta}, \mathcal{M}_k) p(\mathbf{x} | \boldsymbol{\theta}, \mathcal{M}) p(\boldsymbol{\theta} | \mathcal{M}) d\mathbf{x} d\boldsymbol{\theta}$$

- ▶ Prediction with a new data \mathbf{z}

$$p(\mathbf{z} | \mathcal{M}) = \int \int p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}, \mathcal{M}) p(\mathbf{x} | \boldsymbol{\theta}, \mathcal{M}) p(\boldsymbol{\theta} | \mathcal{M}) d\mathbf{x} d\boldsymbol{\theta}$$

Lower Bounding the Marginal Likelihood

Jensen's inequality:

$$\begin{aligned}\ln p(\mathbf{d}|\mathcal{M}_k) &= \ln \int \int p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta}|\mathcal{M}_k) d\mathbf{x} d\boldsymbol{\theta} \\ &= \ln \int \int q(\mathbf{x}, \boldsymbol{\theta}) \frac{p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta}|\mathcal{M}_k)}{q(\mathbf{x}, \boldsymbol{\theta})} d\mathbf{x} d\boldsymbol{\theta} \\ &\geq \int \int q(\mathbf{x}, \boldsymbol{\theta}) \ln \frac{p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta}|\mathcal{M}_k)}{q(\mathbf{x}, \boldsymbol{\theta})} d\mathbf{x} d\boldsymbol{\theta}\end{aligned}$$

Using a factorised approximation for $q(\mathbf{x}, \boldsymbol{\theta}) = q_1(\mathbf{x})q_2(\boldsymbol{\theta})$:

$$\begin{aligned}\ln p(\mathbf{d}|\mathcal{M}_k) &\geq \int \int q_1(\mathbf{x})q_2(\boldsymbol{\theta}) \ln \frac{p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta}|\mathcal{M}_k)}{q_1(\mathbf{x})q_2(\boldsymbol{\theta})} d\mathbf{x} d\boldsymbol{\theta} \\ &= \mathcal{F}_{\mathcal{M}_k}(q_1(\mathbf{x}), q_2(\boldsymbol{\theta}), \mathbf{d})\end{aligned}$$

Maximising this free energy leads to VBA.

Variational Bayesian Learning

$$\begin{aligned}\mathcal{F}_{\mathcal{M}}(q_1(\mathbf{x}), q_2(\boldsymbol{\theta}), \mathbf{d}) &= \int \int q_1(\mathbf{x}) q_2(\boldsymbol{\theta}) \ln \frac{p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta} | \mathcal{M})}{q_1(\mathbf{x}) q_2(\boldsymbol{\theta})} d\mathbf{x} d\boldsymbol{\theta} \\ &= \mathcal{H}(q_1) + \mathcal{H}(q_2) + \langle \ln p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta} | \mathcal{M}) \rangle_{q_1 q_2}\end{aligned}$$

Minimising this lower bound with respect to q_1 and then q_2 leads to EM-like iterative update

$$q_1^{(t+1)}(\mathbf{x}) \propto \exp \left[\langle \ln p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta} | \mathcal{M}) \rangle_{q_2^{(t)}(\boldsymbol{\theta})} \right] \quad \text{E-like step}$$

$$q_2^{(t+1)}(\boldsymbol{\theta}) \propto \exp \left[\langle \ln p(\mathbf{d}, \mathbf{x}, \boldsymbol{\theta} | \mathcal{M}) \rangle_{q_1^{(t+1)}(\mathbf{x})} \right] \quad \text{M-like step}$$

which can also be written as:

$$q_1^{(t+1)}(\mathbf{x}) \propto \exp \left[\langle \ln p(\mathbf{d}, \mathbf{x} | \boldsymbol{\theta}, \mathcal{M}) \rangle_{q_2^{(t)}(\boldsymbol{\theta})} \right] \quad \text{E-like step}$$

$$q_2^{(t+1)}(\boldsymbol{\theta}) \propto p(\boldsymbol{\theta} | \mathcal{M}) \exp \left[\langle \ln p(\mathbf{d}, \mathbf{x} | \boldsymbol{\theta}, \mathcal{M}) \rangle_{q_1^{(t+1)}(\mathbf{x})} \right] \quad \text{M-like step}$$

EM and VBEM algorithms

EM for Marginal MAP estimation

Goal: maximize $p(\theta|\mathbf{d}, \mathcal{M})$ w.r.t. θ

E Step: Compute

$$q_1^{(t+1)}(\mathbf{x}) = p(\mathbf{x}|\mathbf{d}, \theta^{(t)}) \text{ and}$$

$$Q(\theta) = \langle \ln p(\mathbf{d}, \mathbf{x}, \theta | \mathcal{M}) \rangle_{q_1^{(t+1)}(\mathbf{x})}$$

M Step: Maximize

$$\theta^{(t+1)} = \arg \max_{\theta} \{ Q(\theta) \}$$

Variational Bayesian EM

Goal: lower bound $p(\mathbf{d}|\mathcal{M})$

VB-E Step: Compute

$$q_1^{(t+1)}(\mathbf{x}) = p(\mathbf{x}|\mathbf{d}, \phi^{(t)}) \text{ and}$$

$$Q(\theta) = \langle \ln p(\mathbf{d}, \mathbf{x}, \theta | \mathcal{M}) \rangle_{q_1^{(t+1)}(\mathbf{x})}$$

M Step: Maximize

$$q_2^{(t+1)}(\theta) = \exp [Q(\theta)]$$

Properties:

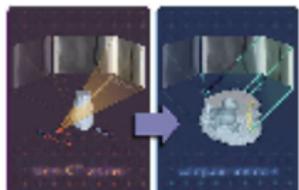
- ▶ VB-EM reduces to EM if $q_2(\theta) = \delta(\theta - \tilde{\theta})$
- ▶ VB-EM has the same complexity than EM
- ▶ If we choose $q_2(\theta)$ in the conjugate family of $p(\mathbf{d}, \mathbf{x}|\theta)$, then ϕ becomes the expected natural parameters
- ▶ The main computational part of both methods is in the E-step. We can use belief propagation, Kalman filter, etc. to do it. In VB-EM, ϕ replaces θ .

Computed Tomography: Seeing inside of a body

- ▶ $f(x, y)$ a section of a real 3D body $f(x, y, z)$
- ▶ $g_\phi(r)$ a line of observed radiography $g_\phi(r, z)$



- ▶ Forward model:
Line integrals or Radon Transform



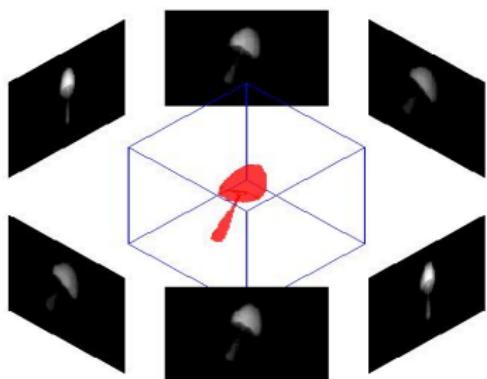
$$\begin{aligned} g_\phi(r) &= \int_{L_{r,\phi}} f(x, y) \, dl + \epsilon_\phi(r) \\ &= \iint f(x, y) \delta(r - x \cos \phi - y \sin \phi) \, dx \, dy + \epsilon_\phi(r) \end{aligned}$$

- ▶ Inverse problem: Image reconstruction

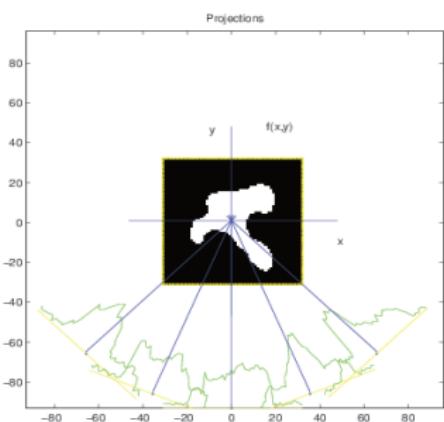
Given the forward model \mathcal{H} (Radon Transform) and
a set of data $g_{\phi_i}(r), i = 1, \dots, M$
find $f(x, y)$

2D and 3D Computed Tomography

3D



2D

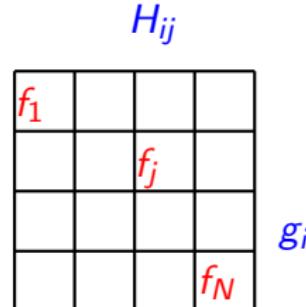
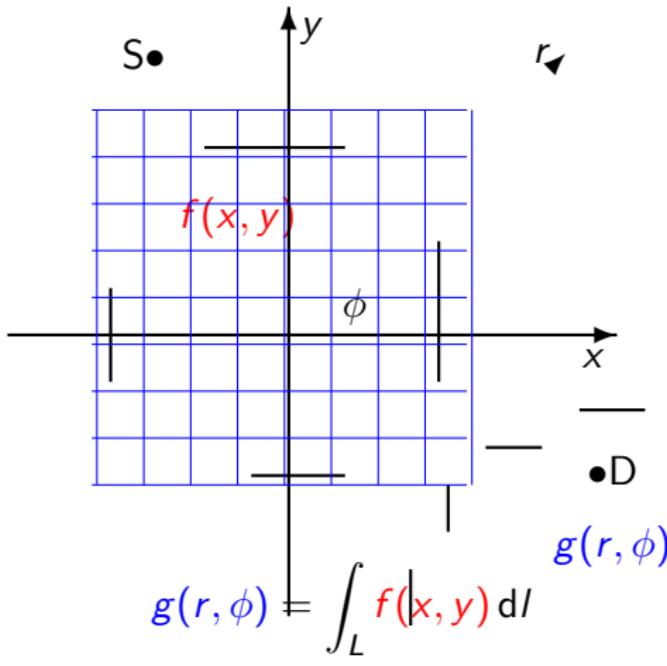


$$g_\phi(r_1, r_2) = \int_{\mathcal{L}_{r_1, r_2, \phi}} f(x, y, z) \, dI \quad g_\phi(r) = \int_{\mathcal{L}_{r, \phi}} f(x, y) \, dI$$

Forward problem: $f(x, y)$ or $f(x, y, z)$ \rightarrow $g_\phi(r)$ or $g_\phi(r_1, r_2)$

Inverse problem: $g_\phi(r)$ or $g_\phi(r_1, r_2)$ \rightarrow $f(x, y)$ or $f(x, y, z)$

Algebraic methods: Discretization



$$f(x, y) = \sum_j f_j b_j(x, y)$$
$$b_j(x, y) = \begin{cases} 1 & \text{if } (x, y) \in \text{ pixel } j \\ 0 & \text{else} \end{cases}$$

$$g_i = \sum_{j=1}^N H_{ij} f_j + \epsilon_i \rightarrow \mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}$$

- \mathbf{H} is huge dimensional: 2D: $10^6 \times 10^6$, 3D: $10^9 \times 10^9$.
- $\mathbf{H}\mathbf{f}$ corresponds to forward projection
- $\mathbf{H}^t\mathbf{g}$ corresponds to Back projection (BP)

Bayesian approach for linear inverse problems

$$\mathcal{M} : \quad \mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}$$

- ▶ Observation model \mathcal{M} + Information on the noise $\boldsymbol{\epsilon}$:

$$p(\mathbf{g}|\mathbf{f}, \theta_1; \mathcal{M}) = p_{\boldsymbol{\epsilon}}(\mathbf{g} - \mathbf{H}\mathbf{f}|\theta_1)$$

- ▶ A priori information $p(\mathbf{f}|\theta_2; \mathcal{M})$

- ▶ Basic Bayes :

$$p(\mathbf{f}|\mathbf{g}, \theta_1, \theta_2; \mathcal{M}) = \frac{p(\mathbf{g}|\mathbf{f}, \theta_1; \mathcal{M}) p(\mathbf{f}|\theta_2; \mathcal{M})}{p(\mathbf{g}|\theta_1, \theta_2; \mathcal{M})}$$

- ▶ Unsupervised:

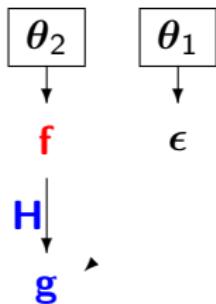
$$p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}, \boldsymbol{\alpha}_0) = \frac{p(\mathbf{g}|\mathbf{f}, \theta_1) p(\mathbf{f}|\theta_2) p(\boldsymbol{\theta}|\boldsymbol{\alpha}_0)}{p(\mathbf{g}|\boldsymbol{\alpha}_0)}, \quad \boldsymbol{\theta} = (\theta_1, \theta_2)$$

- ▶ Hierarchical prior models:

$$p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}|\mathbf{g}, \boldsymbol{\alpha}_0) = \frac{p(\mathbf{g}|\mathbf{f}, \theta_1) p(\mathbf{f}|\mathbf{z}, \theta_2) p(\mathbf{z}|\theta_3) p(\boldsymbol{\theta}|\boldsymbol{\alpha}_0)}{p(\mathbf{g}|\boldsymbol{\alpha}_0)}, \quad \boldsymbol{\theta} = (\theta_1, \theta_2, \theta_3)$$

Bayesian inference for inverse problems

Simple case:



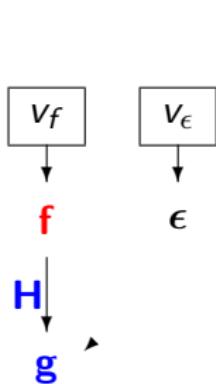
$$\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}$$

$$p(\mathbf{f}|\mathbf{g}, \theta) \propto p(\mathbf{g}|\mathbf{f}, \theta_1) p(\mathbf{f}|\theta_2)$$

– Objective: Infer \mathbf{f}

– MAP: $\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g}, \theta)\}$

– Posterior Mean (PM): $\hat{\mathbf{f}} = \int \mathbf{f} p(\mathbf{f}|\mathbf{g}, \theta) d\mathbf{f}$



Example: Gaussian case:

$$\begin{cases} p(\mathbf{g}|\mathbf{f}, v_\epsilon) = \mathcal{N}(\mathbf{g}|\mathbf{H}\mathbf{f}, v_\epsilon \mathbf{I}) \\ p(\mathbf{f}|v_f) = \mathcal{N}(\mathbf{f}|0, v_f \mathbf{I}) \end{cases} \rightarrow p(\mathbf{f}|\mathbf{g}, \theta) = \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \hat{\Sigma})$$

– MAP: $\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\}$ with
 $J(\mathbf{f}) = \frac{1}{v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \frac{1}{v_f} \|\mathbf{f}\|^2$

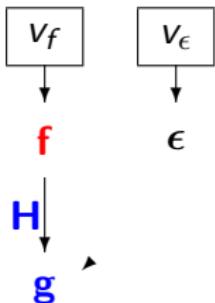
– Posterior Mean (PM)=MAP:

$$\begin{cases} \hat{\mathbf{f}} = (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^t \mathbf{g} \\ \hat{\Sigma} = (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{I})^{-1} \end{cases} \text{ with } \lambda = \frac{v_\epsilon}{v_f}.$$

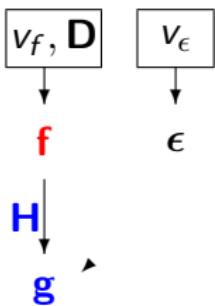
Gaussian model: Simple separable and Markovian

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}$$

Separable Gaussian



Gauss-Markov



$$\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}$$

$$\begin{cases} p(\mathbf{g}|\mathbf{f}, \theta_1) = \mathcal{N}(\mathbf{g}|\mathbf{H}\mathbf{f}, v_\epsilon \mathbf{I}) \\ p(\mathbf{f}|v_f) = \mathcal{N}(\mathbf{f}|0, v_f \mathbf{I}) \end{cases} \xrightarrow{} p(\mathbf{f}|\mathbf{g}, \theta) = \mathcal{N}(\mathbf{f}|\hat{\mathbf{f}}, \hat{\Sigma})$$

– MAP: $\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\}$ with

$$J(\mathbf{f}) = \frac{1}{v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \frac{1}{v_f} \|\mathbf{f}\|^2$$

– Posterior Mean (PM)=MAP:

$$\begin{cases} \hat{\mathbf{f}} = (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^t \mathbf{g} \\ \hat{\Sigma} = v_\epsilon (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{I})^{-1} \end{cases} \text{ with } \lambda = \frac{v_\epsilon}{v_f}.$$

Markovian case:

$$p(\mathbf{f}|v_f, \mathbf{D}) = \mathcal{N}(\mathbf{f}|0, v_f (\mathbf{D}\mathbf{D}^t)^{-1})$$

$$- \text{MAP: } J(\mathbf{f}) = \frac{1}{v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \frac{1}{v_f} \|\mathbf{D}\mathbf{f}\|^2$$

– Posterior Mean (PM)=MAP:

$$\begin{cases} \hat{\mathbf{f}} = (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{D}^t \mathbf{D})^{-1} \mathbf{H}^t \mathbf{g} \\ \hat{\Sigma} = v_\epsilon (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{D}^t \mathbf{D})^{-1} \end{cases} \text{ with } \lambda = \frac{v_\epsilon}{v_f}.$$

Bayesian inference (Unsupervised case)

Unsupervised case: Hyper parameter estimation

$$p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}) \propto p(\mathbf{g} | \mathbf{f}, \boldsymbol{\theta}_1) p(\mathbf{f} | \boldsymbol{\theta}_2) p(\boldsymbol{\theta})$$

– Objective: Infer $(\mathbf{f}, \boldsymbol{\theta})$

$$\text{JMAP: } (\hat{\mathbf{f}}, \hat{\boldsymbol{\theta}}) = \arg \max_{(\mathbf{f}, \boldsymbol{\theta})} \{p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g})\}$$

– Marginalization 1:

$$p(\mathbf{f} | \mathbf{g}) = \int p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}) d\boldsymbol{\theta}$$

– Marginalization 2:

$$p(\boldsymbol{\theta} | \mathbf{g}) = \int p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}) d\mathbf{f} \text{ followed by:}$$

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \{p(\boldsymbol{\theta} | \mathbf{g})\} \rightarrow \hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \left\{ p(\mathbf{f} | \mathbf{g}, \hat{\boldsymbol{\theta}}) \right\}$$

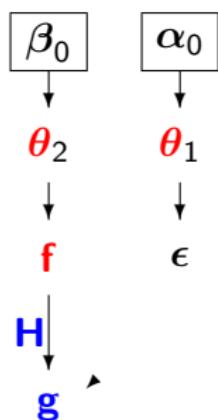
– MCMC Gibbs sampling:

$$\mathbf{f} \sim p(\mathbf{f} | \boldsymbol{\theta}, \mathbf{g}) \rightarrow \boldsymbol{\theta} \sim p(\boldsymbol{\theta} | \mathbf{f}, \mathbf{g}) \text{ until convergence}$$

Use samples generated to compute mean and variances

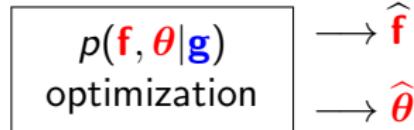
– VBA: Approximate $p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g})$ by $q_1(\mathbf{f}) q_2(\boldsymbol{\theta})$

Use $q_1(\mathbf{f})$ to infer \mathbf{f} and $q_2(\boldsymbol{\theta})$ to infer $\boldsymbol{\theta}$

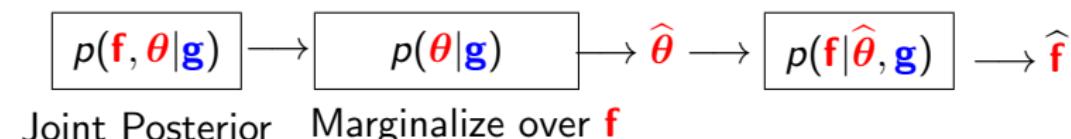


JMAP, Marginalization, VBA

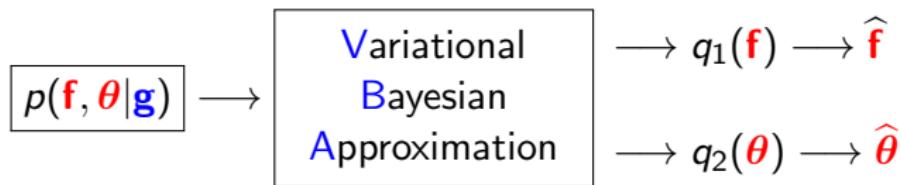
- ▶ JMAP:



- ▶ Marginalization



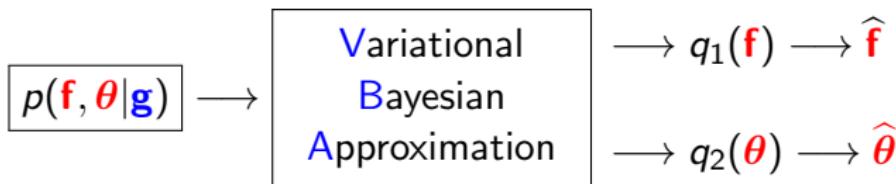
- ▶ Variational Bayesian Approximation



Variational Bayesian Approximation

- ▶ Approximate $p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g})$ by $q(\mathbf{f}, \boldsymbol{\theta}) = q_1(\mathbf{f}) q_2(\boldsymbol{\theta})$ and then use them for any inferences on \mathbf{f} and $\boldsymbol{\theta}$ respectively.
- ▶ Criterion $\text{KL}(q(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}) : p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}))$
$$\text{KL}(q : p) = \int \int q \ln \frac{q}{p} = \int \int q_1 q_2 \ln \frac{q_1 q_2}{p}$$
- ▶ Iterative algorithm $q_1 \rightarrow q_2 \rightarrow q_1 \rightarrow q_2, \dots$

$$\begin{cases} \hat{q}_1(\mathbf{f}) & \propto \exp \left[\langle \ln p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta}; \mathcal{M}) \rangle_{\hat{q}_2(\boldsymbol{\theta})} \right] \\ \hat{q}_2(\boldsymbol{\theta}) & \propto \exp \left[\langle \ln p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta}; \mathcal{M}) \rangle_{\hat{q}_1(\mathbf{f})} \right] \end{cases}$$



Variational Bayesian Approximation

$$p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta} | \mathcal{M}) = p(\mathbf{g} | \mathbf{f}, \boldsymbol{\theta}, \mathcal{M}) p(\mathbf{f} | \boldsymbol{\theta}, \mathcal{M}) p(\boldsymbol{\theta} | \mathcal{M})$$

$$p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}, \mathcal{M}) = \frac{p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta} | \mathcal{M})}{p(\mathbf{g} | \mathcal{M})}$$

$$\text{KL}(q : p) = \iint q(\mathbf{f}, \boldsymbol{\theta}) \ln \frac{p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}; \mathcal{M})}{q(\mathbf{f}, \boldsymbol{\theta})} d\mathbf{f} d\boldsymbol{\theta}$$

$$\begin{aligned} p(\mathbf{g} | \mathcal{M}) &= \iint q(\mathbf{f}, \boldsymbol{\theta}) \frac{p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta} | \mathcal{M})}{q(\mathbf{f}, \boldsymbol{\theta})} d\mathbf{f} d\boldsymbol{\theta} \\ &\geq \iint q(\mathbf{f}, \boldsymbol{\theta}) \ln \frac{p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta} | \mathcal{M})}{q(\mathbf{f}, \boldsymbol{\theta})} d\mathbf{f} d\boldsymbol{\theta} \end{aligned}$$

Free energy:

$$\mathcal{F}(q) = \iint q(\mathbf{f}, \boldsymbol{\theta}) \ln \frac{p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta} | \mathcal{M})}{q(\mathbf{f}, \boldsymbol{\theta})} d\mathbf{f} d\boldsymbol{\theta}$$

Evidence of the model \mathcal{M} :

$$p(\mathbf{g} | \mathcal{M}) = \mathcal{F}(q) + \text{KL}(q : p)$$

VBA: Separable Approximation

$$p(\mathbf{g}|\mathcal{M}) = \mathcal{F}(q) + \text{KL}(q : p)$$

$$q(\mathbf{f}, \boldsymbol{\theta}) = q_1(\mathbf{f}) q_2(\boldsymbol{\theta})$$

Minimizing $\text{KL}(q : p) = \text{Maximizing } \mathcal{F}(q)$

$$(\hat{q}_1, \hat{q}_2) = \arg \min_{(q_1, q_2)} \{\text{KL}(q_1 q_2 : p)\} = \arg \max_{(q_1, q_2)} \{\mathcal{F}(q_1 q_2)\}$$

$\text{KL}(q_1 q_2 : p)$ is convex wrt q_1 when q_2 is fixed and vice versa:

$$\begin{cases} \hat{q}_1 = \arg \min_{q_1} \{\text{KL}(q_1 \hat{q}_2 : p)\} = \arg \max_{q_1} \{\mathcal{F}(q_1 \hat{q}_2)\} \\ \hat{q}_2 = \arg \min_{q_2} \{\text{KL}(\hat{q}_1 q_2 : p)\} = \arg \max_{q_2} \{\mathcal{F}(\hat{q}_1 q_2)\} \end{cases}$$

$$\begin{cases} \hat{q}_1(\mathbf{f}) \propto \exp \left[\langle \ln p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta}; \mathcal{M}) \rangle_{\hat{q}_2(\boldsymbol{\theta})} \right] \\ \hat{q}_2(\boldsymbol{\theta}) \propto \exp \left[\langle \ln p(\mathbf{g}, \mathbf{f}, \boldsymbol{\theta}; \mathcal{M}) \rangle_{\hat{q}_1(\mathbf{f})} \right] \end{cases}$$

BVA: Choice of family of laws q_1 and q_2

- Case 1 : \rightarrow Joint MAP

$$\begin{cases} \hat{q}_1(\mathbf{f}|\tilde{\mathbf{f}}) = \delta(\mathbf{f} - \tilde{\mathbf{f}}) \\ \hat{q}_2(\boldsymbol{\theta}|\tilde{\boldsymbol{\theta}}) = \delta(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) \end{cases} \rightarrow \begin{cases} \tilde{\mathbf{f}} = \arg \max_{\mathbf{f}} \left\{ p(\mathbf{f}, r\tilde{\boldsymbol{\theta}} | \mathbf{g}; \mathcal{M}) \right\} \\ \tilde{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \left\{ p(\tilde{\mathbf{f}}, \boldsymbol{\theta} | \mathbf{g}; \mathcal{M}) \right\} \end{cases}$$

- Case 2 : \rightarrow EM

$$\begin{cases} \hat{q}_1(\mathbf{f}) \propto p(\mathbf{f} | \boldsymbol{\theta}, \mathbf{g}) \\ \hat{q}_2(\boldsymbol{\theta} | \tilde{\boldsymbol{\theta}}) = \delta(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}) \end{cases} \rightarrow \begin{cases} Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) = \langle \ln p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}; \mathcal{M}) \rangle_{q_1(\mathbf{f} | \tilde{\boldsymbol{\theta}})} \\ \tilde{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \left\{ Q(\boldsymbol{\theta}, \tilde{\boldsymbol{\theta}}) \right\} \end{cases}$$

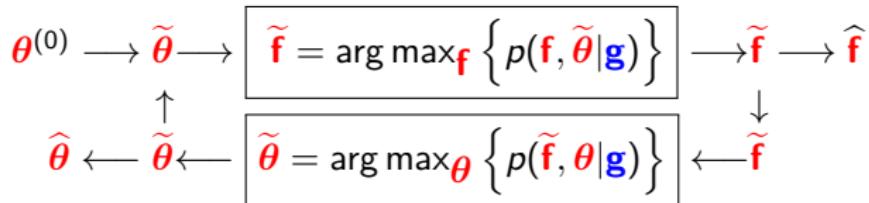
- Appropriate choice for inverse problems

$$\begin{cases} \hat{q}_1(\mathbf{f}) \propto p(\mathbf{f} | \tilde{\boldsymbol{\theta}}, \mathbf{g}; \mathcal{M}) \\ \hat{q}_2(\boldsymbol{\theta}) \propto p(\boldsymbol{\theta} | \mathbf{f}, \mathbf{g}; \mathcal{M}) \end{cases} \rightarrow \begin{cases} \text{Accounts for the uncertainties of} \\ \hat{\boldsymbol{\theta}} \text{ for } \hat{\mathbf{f}} \text{ and vice versa.} \end{cases}$$

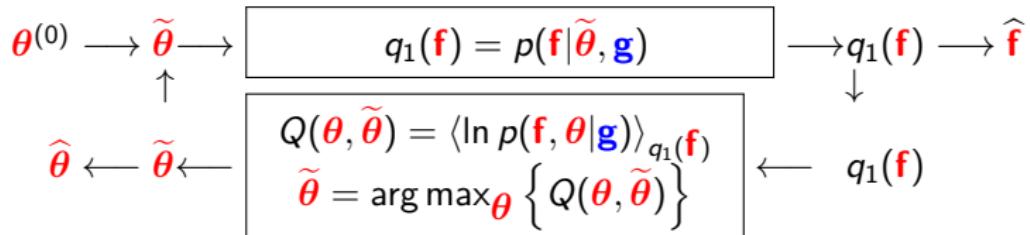
Exponential families, Conjugate priors

JMAP, EM and VBA

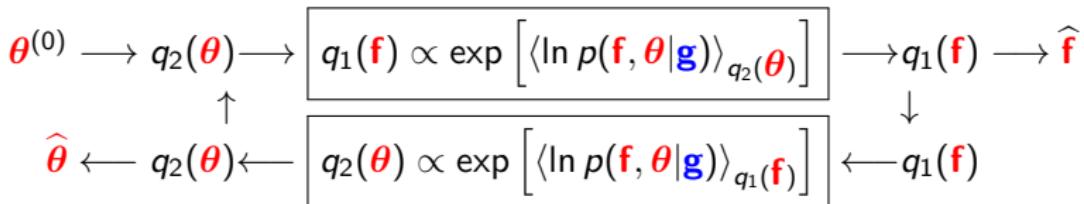
JMAP Alternate optimization Algorithm:



EM:



VBA:



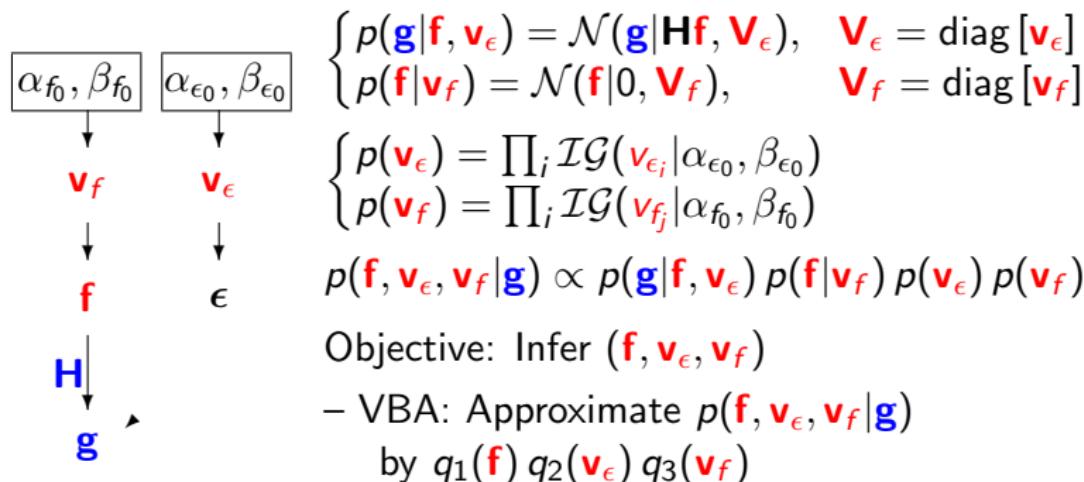
Non stationary noise and sparsity enforcing model

- Non stationary noise:

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}, \quad \epsilon_i \sim \mathcal{N}(\epsilon_i | 0, v_{\epsilon_i}) \rightarrow \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon} | 0, \mathbf{V}_{\epsilon} = \text{diag} [v_{\epsilon 1}, \dots, v_{\epsilon M}])$$

- Student-t prior model and its equivalent IGSM :

$$f_j | v_{f_j} \sim \mathcal{N}(f_j | 0, v_{f_j}) \text{ and } v_{f_j} \sim \mathcal{IG}(v_{f_0} | \alpha_{f_0}, \beta_{f_0}) \rightarrow f_j \sim \mathcal{St}(f_j | \alpha_{f_0}, \beta_{f_0})$$



Sparse model in a Transform domain 1

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}, \quad \mathbf{f} = \mathbf{D}\mathbf{z}, \quad \mathbf{z} \text{ sparse}$$

$$\begin{cases} p(\mathbf{g}|\mathbf{z}, \mathbf{v}_\epsilon) = \mathcal{N}(\mathbf{g}|\mathbf{H}\mathbf{D}\mathbf{f}, \mathbf{V}_\epsilon \mathbf{I}) \\ p(\mathbf{z}|\mathbf{v}_z) = \mathcal{N}(\mathbf{z}|0, \mathbf{V}_z), \quad \mathbf{V}_z = \text{diag}[\mathbf{v}_z] \end{cases}$$

$$\alpha_{z_0}, \beta_{z_0}$$

$$\begin{cases} p(\mathbf{v}_\epsilon) = \mathcal{IG}(\mathbf{v}_\epsilon | \alpha_{\epsilon_0}, \beta_{\epsilon_0}) \\ p(\mathbf{v}_z) = \prod_i \mathcal{IG}(\mathbf{v}_{z_j} | \alpha_{z_0}, \beta_{z_0}) \end{cases}$$

$$p(\mathbf{z}, \mathbf{v}_\epsilon, \mathbf{v}_z, \mathbf{v}_\xi | \mathbf{g}) \propto p(\mathbf{g}|\mathbf{z}, \mathbf{v}_\epsilon) p(\mathbf{z}|\mathbf{v}_z) p(\mathbf{v}_\epsilon) p(\mathbf{v}_z) p(\mathbf{v}_\xi)$$

- JMAP:

$$(\hat{\mathbf{z}}, \hat{\mathbf{v}}_\epsilon, \hat{\mathbf{v}}_z) = \arg \max_{(\mathbf{z}, \mathbf{v}_\epsilon, \mathbf{v}_z)} \{p(\mathbf{z}, \mathbf{v}_\epsilon, \mathbf{v}_z | \mathbf{g})\}$$

\mathbf{D}

\downarrow

Alternate optimization:

\mathbf{f}

ϵ

$\hat{\mathbf{z}} = \arg \min_{\mathbf{z}} \{J(\mathbf{z})\}$ with:

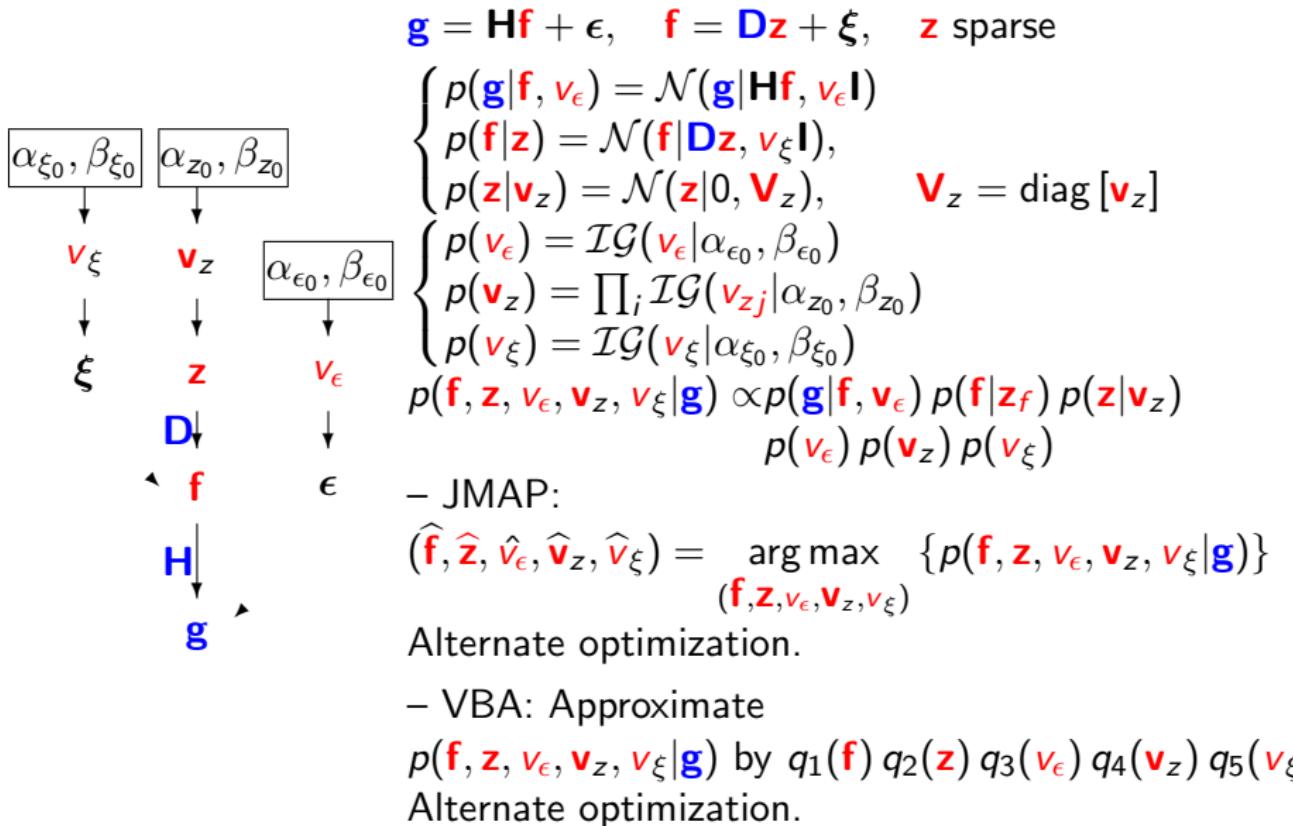
$$\begin{cases} J(\mathbf{z}) = \frac{1}{2\hat{\mathbf{v}}_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{D}\mathbf{z}\|^2 + \|\mathbf{V}_z^{-1/2}\mathbf{z}\|^2 \\ \hat{\mathbf{v}}_{z_j} = \frac{\beta_{z_0} + \hat{z}_j^2}{\alpha_{z_0} + 1/2} \\ \hat{\mathbf{v}}_\epsilon = \frac{\beta_{\epsilon_0} + \|\mathbf{g} - \mathbf{H}\mathbf{D}\hat{\mathbf{z}}\|^2}{\alpha_{\epsilon_0} + M/2} \end{cases}$$

- VBA: Approximate

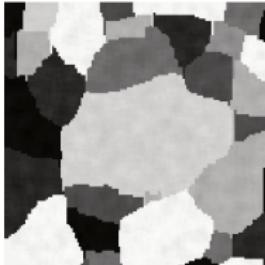
$$p(\mathbf{z}, \mathbf{v}_\epsilon, \mathbf{v}_z, \mathbf{v}_\xi | \mathbf{g}) \text{ by } q_1(\mathbf{z}) q_2(\mathbf{v}_\epsilon) q_3(\mathbf{v}_z)$$

Alternate optimization.

Sparse model in a Transform domain 2



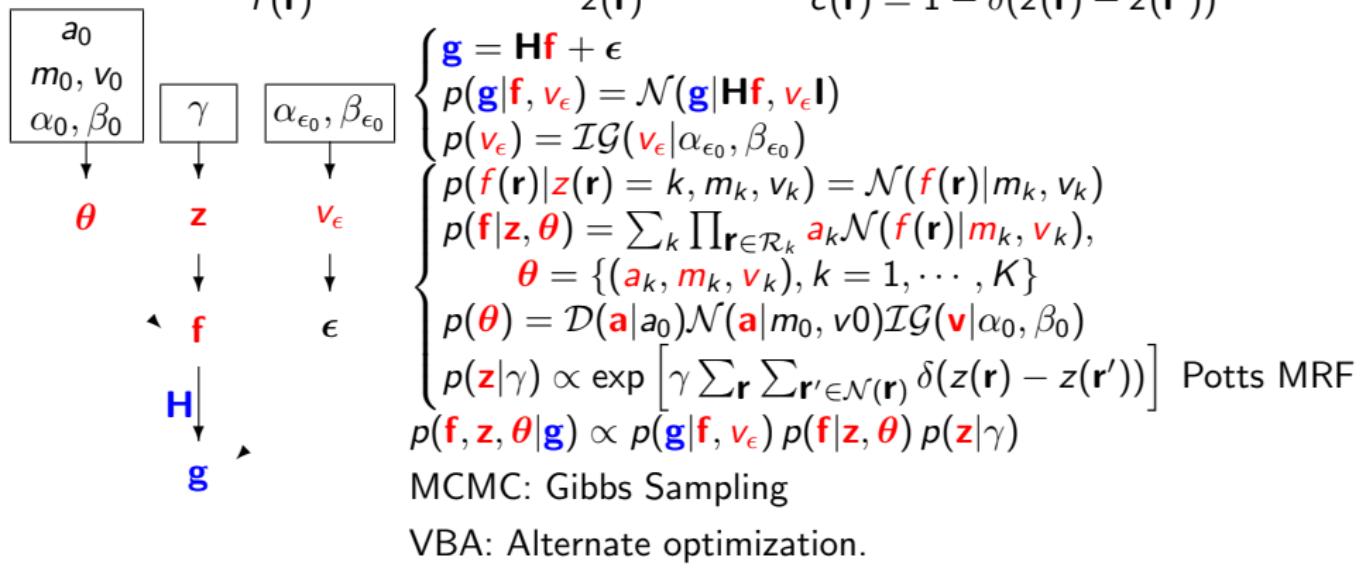
Gauss-Markov-Potts prior models for images



$f(\mathbf{r})$

$z(\mathbf{r})$

$$c(\mathbf{r}) = 1 - \delta(z(\mathbf{r}) - z(\mathbf{r}'))$$



Mixture Models

1. Mixture models
2. Different problems related to classification and clustering
 - ▶ Training
 - ▶ Supervised classification
 - ▶ Semi-supervised classification
 - ▶ Clustering or unsupervised classification
3. Mixture of Gaussian (MoG)
4. Mixture of Student-t (MoSt)
5. Variational Bayesian Approximation (VBA)
6. VBA for Mixture of Gaussian
7. VBA for Mixture of Student-t
8. Conclusion

Mixture models

- ▶ General mixture model

$$p(\mathbf{x}|\mathbf{a}, \Theta, K) = \sum_{k=1}^K a_k p_k(\mathbf{x}_k|\theta_k), \quad 0 < a_k < 1, \quad \sum_{k=1}^K a_k = 1$$

- ▶ Same family $p_k(\mathbf{x}_k|\theta_k) = p(\mathbf{x}_k|\theta_k)$, $\forall k$
- ▶ Gaussian $p(\mathbf{x}_k|\theta_k) = \mathcal{N}(\mathbf{x}_k|\mu_k, \mathbf{V}_k)$ with $\theta_k = (\mu_k, \mathbf{V}_k)$
- ▶ Data $\mathbf{X} = \{\mathbf{x}_n, n = 1, \dots, N\}$ where each element \mathbf{x}_n can be in one of the K classes c_n .
- ▶ $a_k = p(c_n = k)$, $\mathbf{a} = \{a_k, k = 1, \dots, K\}$,
 $\Theta = \{\theta_k, k = 1, \dots, K\}$, $\mathbf{c} = \{c_n, n = 1, \dots, N\}$

$$p(\mathbf{X}, \mathbf{c}|\mathbf{a}, \Theta) = \prod_{n=1}^N p(\mathbf{x}_n, c_n = k | a_k, \theta_k)$$

Different problems

- ▶ **Training:**

Given a set of (training) data \mathbf{X} and classes \mathbf{c} , estimate the parameters \mathbf{a} and Θ .

- ▶ **Supervised classification:**

Given a sample \mathbf{x}_m and the parameters K , \mathbf{a} and Θ determine its class

$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, \mathbf{a}, \Theta, K)\}.$$

- ▶ **Semi-supervised classification (Proportions are not known):**

Given sample \mathbf{x}_m and the parameters K and Θ , determine its class

$$k^* = \arg \max_k \{p(c_m = k | \mathbf{x}_m, \Theta, K)\}.$$

- ▶ **Clustering or unsupervised classification (Number of classes K is not known):**

Given a set of data \mathbf{X} , determine K and \mathbf{c} .

Training

- ▶ Given a set of (training) data \mathbf{X} and classes \mathbf{c} , estimate the parameters \mathbf{a} and $\boldsymbol{\Theta}$.
- ▶ Maximum Likelihood (ML):

$$(\hat{\mathbf{a}}, \hat{\boldsymbol{\Theta}}) = \arg \max_{(\mathbf{a}, \boldsymbol{\Theta})} \{ p(\mathbf{X}, \mathbf{c} | \mathbf{a}, \boldsymbol{\Theta}, K) \}.$$

- ▶ Bayesian: Assign priors $p(\mathbf{a}|K)$ and $p(\boldsymbol{\Theta}|K) = \prod_{k=1}^K p(\theta_k|K)$ and write the expression of the joint posterior laws:

$$p(\mathbf{a}, \boldsymbol{\Theta} | \mathbf{X}, \mathbf{c}, K) = \frac{p(\mathbf{X}, \mathbf{c} | \mathbf{a}, \boldsymbol{\Theta}, K) p(\mathbf{a}|K) p(\boldsymbol{\Theta}|K)}{p(\mathbf{X}, \mathbf{c}|K)}$$

where

$$p(\mathbf{X}, \mathbf{c}|K) = \iint p(\mathbf{X}, \mathbf{c} | \mathbf{a}, \boldsymbol{\Theta} | K) p(\mathbf{a}|K) p(\boldsymbol{\Theta}|K) d\mathbf{a} d\boldsymbol{\Theta}$$

- ▶ Infer on \mathbf{a} and $\boldsymbol{\Theta}$ either as the Maximum A Posteriori (MAP) or Posterior Mean (PM).

Supervised classification

- Given a sample \mathbf{x}_m and the parameters K , \mathbf{a} and Θ determine

$$p(\mathbf{c}_m = k | \mathbf{x}_m, \mathbf{a}, \Theta, K) = \frac{p(\mathbf{x}_m, \mathbf{c}_m = k | \mathbf{a}, \Theta, K)}{p(\mathbf{x}_m | \mathbf{a}, \Theta, K)}$$

where $p(\mathbf{x}_m, \mathbf{c}_m = k | \mathbf{a}, \Theta, K) = a_k p(\mathbf{x}_m | \theta_k)$ and

$$p(\mathbf{x}_m | \mathbf{a}, \Theta, K) = \sum_{k=1}^K a_k p(\mathbf{x}_m | \theta_k)$$

- Best class k^* :

$$k^* = \arg \max_k \{p(\mathbf{c}_m = k | \mathbf{x}_m, \mathbf{a}, \Theta, K)\}$$

Semi-supervised classification

- Given sample \mathbf{x}_m and the parameters K and Θ (not the proportions \mathbf{a}), determine the probabilities

$$p(\mathbf{c}_m = k | \mathbf{x}_m, \Theta, K) = \frac{p(\mathbf{x}_m, \mathbf{c}_m = k | \Theta, K)}{p(\mathbf{x}_m | \Theta, K)}$$

where

$$p(\mathbf{x}_m, \mathbf{c}_m = k | \Theta, K) = \int p(\mathbf{x}_m, \mathbf{c}_m = k | \mathbf{a}, \Theta, K) p(\mathbf{a} | K) d\mathbf{a}$$

and

$$p(\mathbf{x}_m | \Theta, K) = \sum_{k=1}^K p(\mathbf{x}_m, \mathbf{c}_m = k | \Theta, K)$$

- Best class k^* , for example the MAP solution:

$$k^* = \arg \max_k \{p(\mathbf{c}_m = k | \mathbf{x}_m, \Theta, K)\}.$$

Clustering or non-supervised classification

- Given a set of data \mathbf{X} , determine K and \mathbf{c} .
- Determination of the number of classes:

$$p(K = L | \mathbf{X}) = \frac{p(\mathbf{X}, K = L)}{p(\mathbf{X})} = \frac{p(\mathbf{X}|K = L) p(K = L)}{p(\mathbf{X})}$$

and

$$p(\mathbf{X}) = \sum_{L=1}^{L_0} p(K = L) p(\mathbf{X}|K = L),$$

where L_0 is the a priori maximum number of classes and

$$p(\mathbf{X}|K = L) = \int \int \prod_n \prod_{k=1}^L a_k p(\mathbf{x}_n, c_n = k | \boldsymbol{\theta}_k) p(\mathbf{a}|K) p(\boldsymbol{\Theta}|K) d\mathbf{a} d\boldsymbol{\Theta}.$$

- When K and \mathbf{c} are determined, we can also determine the characteristics of those classes \mathbf{a} and $\boldsymbol{\Theta}$.

Mixture of Gaussian and Mixture of Student-t

$$p(\mathbf{x}|\mathbf{a}, \Theta, K) = \sum_{k=1}^K a_k p(\mathbf{x}_k|\boldsymbol{\theta}_k), \quad 0 < a_k < 1, \quad \sum_{k=1}^K a_k = 1$$

- ▶ Mixture of Gaussian (MoG)

$$p(\mathbf{x}_k|\boldsymbol{\theta}_k) = \mathcal{N}(\mathbf{x}_k|\boldsymbol{\mu}_k, \mathbf{V}_k), \quad \boldsymbol{\theta}_k = (\boldsymbol{\mu}_k, \mathbf{V}_k)$$

$$\mathcal{N}(\mathbf{x}_k|\boldsymbol{\mu}_k, \mathbf{V}_k) = (2\pi)^{-\frac{p}{2}} |\mathbf{V}_k|^{-\frac{1}{2}} \exp \left[\frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu}_k)' \mathbf{V}_k^{-1} (\mathbf{x}_k - \boldsymbol{\mu}_k) \right]$$

- ▶ Mixture of Student-t (MoSt)

$$p(\mathbf{x}_k|\boldsymbol{\theta}_k) = \mathcal{T}(\mathbf{x}_k|\nu_k, \boldsymbol{\mu}_k, \mathbf{V}_k), \quad \boldsymbol{\theta}_k = (\nu_k, \boldsymbol{\mu}_k, \mathbf{V}_k)$$

$$\mathcal{T}(\mathbf{x}_k|\nu, \boldsymbol{\mu}_k, \mathbf{V}_k) = \frac{\Gamma \left[\frac{(\nu_k+p)}{2} \right]}{\Gamma \left(\frac{\nu_k}{2} \right) \nu^{\frac{p}{2}} \pi^{\frac{p}{2}}} |\mathbf{V}_k|^{-\frac{1}{2}} \left[1 + \frac{1}{\nu} (\mathbf{x}_k - \boldsymbol{\mu}_k)' \mathbf{V}_k^{-1} (\mathbf{x}_k - \boldsymbol{\mu}_k) \right]^{-\frac{(\nu+p)}{2}}$$

Mixture of Student-t model

- ▶ Student-t and its Infinite Gaussian Scaled Model (IGSM):

$$\mathcal{T}(\mathbf{x}|\nu, \boldsymbol{\mu}, \mathbf{V}) = \int_0^\infty \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, u^{-1}\mathbf{V}) \mathcal{G}(u|\frac{\nu}{2}, \frac{\nu}{2}) dz$$

where

$$\begin{aligned}\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{V}) &= |2\pi\mathbf{V}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \mathbf{V}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \\ &= |2\pi\mathbf{V}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} \text{Tr} \{ (\mathbf{x} - \boldsymbol{\mu}) \mathbf{V}^{-1} (\mathbf{x} - \boldsymbol{\mu})' \} \right]\end{aligned}$$

and

$$\mathcal{G}(u|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} u^{\alpha-1} \exp[-\beta u].$$

- ▶ Mixture of generalized Student-t: $\mathcal{T}(\mathbf{x}|\alpha, \beta, \boldsymbol{\mu}, \mathbf{V})$

$$p(\mathbf{x}|\{a_k, \boldsymbol{\mu}_k, \mathbf{V}_k, \alpha_k, \beta_k\}, K) = \sum_{k=1}^K a_k \mathcal{T}(\mathbf{x}_n|\alpha_k, \beta_k, \boldsymbol{\mu}_k, \mathbf{V}_k).$$

Mixture of Gaussian model

- ▶ Introducing $z_{nk} \in \{0, 1\}$, $\mathbf{z}_k = \{z_{nk}, n = 1, \dots, N\}$,
 $\mathbf{Z} = \{\mathbf{z}_k\}$ with $P(z_{nk} = 1) = P(c_n = k) = a_k$,
 $\boldsymbol{\theta}_k = \{a_k, \boldsymbol{\mu}_k, \mathbf{V}_k\}$, $\boldsymbol{\Theta} = \{\boldsymbol{\theta}_k, k = 1, \dots, K\}$
- ▶ Assigning the priors $p(\boldsymbol{\Theta}) = \prod_k p(\boldsymbol{\theta}_k)$, we can write:

$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K) = \prod_n \sum_k a_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{V}_k) (1 - \delta(z_{nk})) \prod_k p(\boldsymbol{\theta}_k)$$

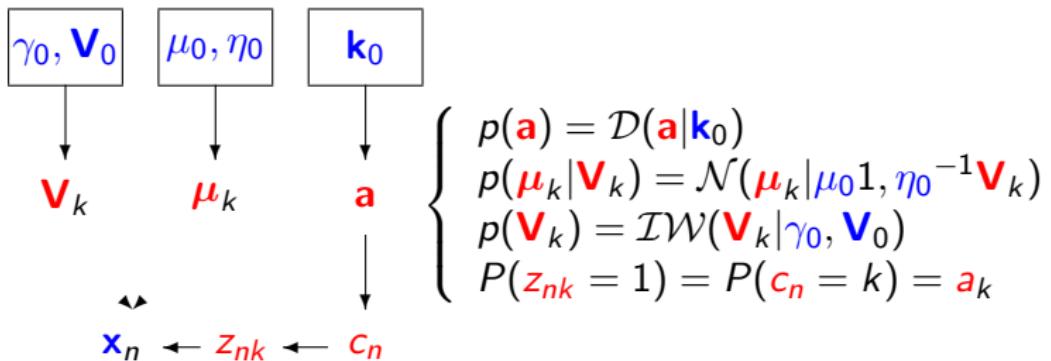
$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K) = \prod_n \prod_k [a_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{V}_k)]^{z_{nk}} p(\boldsymbol{\theta}_k)$$

- ▶ Joint posterior law:

$$p(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | \mathbf{X}, K) = \frac{p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K)}{p(\mathbf{X} | K)}.$$

- ▶ The main task now is to propose some **approximations** to it in such a way that we can use it easily in all the above mentioned tasks of classification or clustering.

Hierarchical graphical model for Mixture of Gaussian



$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \Theta | K) = \prod_n \prod_k [a_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{V}_k)]^{z_{nk}} p(a_k) p(\boldsymbol{\mu}_k | \mathbf{V}_k) p(\mathbf{V}_k)$$

Mixture of Student-t model

- ▶ Introducing $\mathbf{U} = \{u_{nk}\}$

$$\boldsymbol{\theta}_k = \{\alpha_k, \beta_k, \mathbf{a}_k, \boldsymbol{\mu}_k, \mathbf{V}_k\}, \Theta = \{\boldsymbol{\theta}_k, k = 1, \dots, K\}$$

- ▶ Assigning the priors $p(\Theta) = \prod_k p(\boldsymbol{\theta}_k)$, we can write:

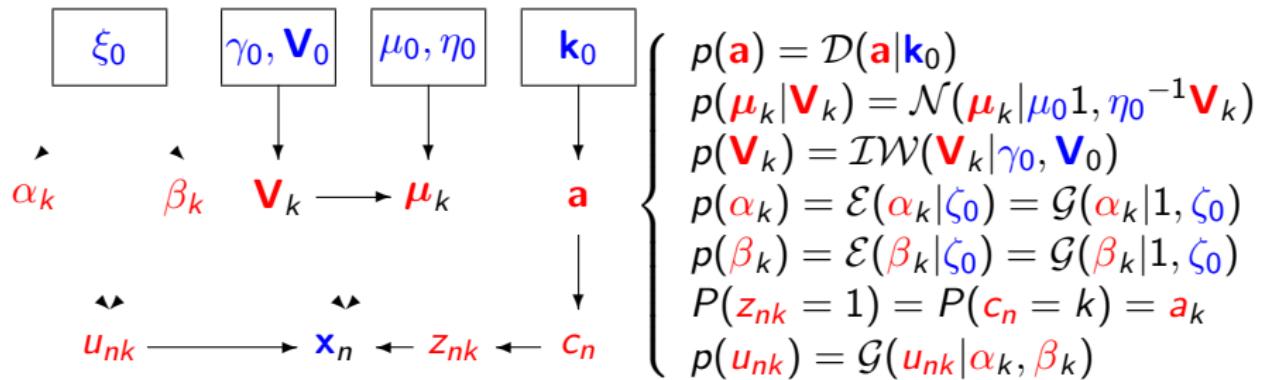
$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \mathbf{U}, \Theta | K) = \prod_n \prod_k [a_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{u}_{n,k}^{-1} \mathbf{V}_k) \mathcal{G}(u_{nk} | \alpha_k, \beta_k)]^{z_{nk}} p(\boldsymbol{\theta}_k)$$

- ▶ Joint posterior law:

$$p(\mathbf{c}, \mathbf{Z}, \mathbf{U}, \Theta | \mathbf{X}, K) = \frac{p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \mathbf{U}, \Theta | K)}{p(\mathbf{X} | K)}.$$

- ▶ The main task now is to propose some **approximations** to it in such a way that we can use it easily in all the above mentioned tasks of classification or clustering.

Hierarchical graphical model for Mixture of Student-t



$$p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \mathbf{U}, \Theta | K) = \prod_n \prod_k [a_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{V}_k) \mathcal{G}(u_{nk} | \alpha_k, \beta_k)]^{z_{nk}} \\ p(\mathbf{a}_k) p(\boldsymbol{\mu}_k | \mathbf{V}_k) p(\mathbf{V}_k) p(\alpha_k) p(\beta_k)$$

Variational Bayesian Approximation (VBA)

- ▶ Main idea: to propose easy computational approximations:
 $q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}) = q(\mathbf{c}, \mathbf{Z})q(\boldsymbol{\Theta})$ for $p(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | \mathbf{X}, K)$ for MoG model,
or
 $q(\mathbf{c}, \mathbf{Z}, \mathbf{U}, \boldsymbol{\Theta}) = q(\mathbf{c}, \mathbf{Z}, \mathbf{U})q(\boldsymbol{\Theta})$ for $p(\mathbf{c}, \mathbf{Z}, \mathbf{U}, \boldsymbol{\Theta} | \mathbf{X}, K)$ for MoSt model.
- ▶ Criterion:

$$\text{KL}(q : p) = -\mathcal{F}(q) + \ln p(\mathbf{X}|K)$$

where

$$\mathcal{F}(q) = \langle -\ln p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K) \rangle_q$$

or

$$\mathcal{F}(q) = \langle -\ln p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \mathbf{U}, \boldsymbol{\Theta} | K) \rangle_q$$

- ▶ Maximizing $\mathcal{F}(q)$ or minimizing $\text{KL}(q : p)$ are equivalent and both give an upper bound to the evidence of the model $\ln p(\mathbf{X}|K)$.
- ▶ When the optimum q^* is obtained, $\mathcal{F}(q^*)$ can be used as a criterion for model selection.

Proposed VBA for Mixture of Student-t priors model

- ▶ Dirichlet

$$\mathcal{D}(\mathbf{a}|\mathbf{k}) = \frac{\Gamma(\sum_I k_I)}{\prod_I \Gamma(k_I)} \prod_I a_I^{k_I - 1}$$

- ▶ Exponential

$$\mathcal{E}(t|\zeta_0) = \zeta_0 \exp[-\zeta_0 t]$$

- ▶ Gamma

$$\mathcal{G}(t|a, b) = \frac{b^a}{\Gamma(a)} t^{a-1} \exp[-bt]$$

- ▶ Inverse Wishart

$$\mathcal{IW}(\mathbf{V}|\gamma, \gamma \boldsymbol{\Delta}) = \frac{\frac{1}{2} |\boldsymbol{\Delta}|^{\gamma/2} \exp\left[-\frac{1}{2} \text{Tr}\{\boldsymbol{\Delta} \mathbf{V}^{-1}\}\right]}{\Gamma_D(\gamma/2) |\mathbf{V}|^{\frac{\gamma+D+1}{2}}}.$$

Expressions of q

$$\begin{aligned} q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta}) &= q(\mathbf{c}, \mathbf{Z}) q(\boldsymbol{\Theta}) \\ &= \prod_n \prod_k [q(c_n = k | z_{nk}) q(z_{nk})] \\ &\quad \prod_k [q(\alpha_k) q(\beta_k) q(\boldsymbol{\mu}_k | \mathbf{V}_k) q(\mathbf{V}_k)] q(\mathbf{a}). \end{aligned}$$

with:

$$\left\{ \begin{array}{l} q(\mathbf{a}) = \mathcal{D}(\mathbf{a} | \tilde{\mathbf{k}}), \quad \tilde{\mathbf{k}} = [\tilde{k}_1, \dots, \tilde{k}_K] \\ q(\alpha_k) = \mathcal{G}(\alpha_k | \tilde{\zeta}_k, \tilde{\eta}_k) \\ q(\beta_k) = \mathcal{G}(\beta_k | \tilde{\zeta}_k, \tilde{\eta}_k) \\ q(\boldsymbol{\mu}_k | \mathbf{V}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \tilde{\boldsymbol{\mu}}, \tilde{\eta}^{-1} \mathbf{V}_k) \\ q(\mathbf{V}_k) = \mathcal{IW}(\mathbf{V}_k | \tilde{\gamma}, \tilde{\Sigma}) \end{array} \right.$$

With these choices, we have

$$\mathcal{F}(q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta})) = \langle \ln p(\mathbf{X}, \mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta} | K) \rangle_{q(\mathbf{c}, \mathbf{Z}, \boldsymbol{\Theta})} = \prod_k \prod_n \mathcal{F}_{1_{kn}} + \prod_k \mathcal{F}_{2_k}$$

$$\mathcal{F}_{1_{kn}} = \langle \ln p(\mathbf{x}_n, c_n, z_{nk}, \boldsymbol{\theta}_k) \rangle_{q(c_n=k|z_{nk})q(z_{nk})}$$

VBA Algorithm step

Expressions of the updating expressions of the tilded parameters are obtained by following three steps:

- ▶ **E step:** Optimizing \mathcal{F} with respect to $q(\mathbf{c}, \mathbf{Z})$ when keeping $q(\boldsymbol{\Theta})$ fixed, we obtain the expression of $q(c_n = k | z_{nk}) = \tilde{a}_k$, $q(z_{nk}) = \mathcal{G}(z_{nk} | \tilde{\alpha}_k, \tilde{\beta}_k)$.
- ▶ **M step:** Optimizing \mathcal{F} with respect to $q(\boldsymbol{\Theta})$ when keeping $q(\mathbf{c}, \mathbf{Z})$ fixed, we obtain the expression of
 $q(\mathbf{a}) = \mathcal{D}(\mathbf{a} | \tilde{\mathbf{k}})$, $\tilde{\mathbf{k}} = [\tilde{k}_1, \dots, \tilde{k}_K]$, $q(\alpha_k) = \mathcal{G}(\alpha_k | \tilde{\zeta}_k, \tilde{\eta}_k)$,
 $q(\beta_k) = \mathcal{G}(\beta_k | \tilde{\zeta}_k, \tilde{\eta}_k)$, $q(\boldsymbol{\mu}_k | \mathbf{V}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \tilde{\boldsymbol{\mu}}, \tilde{\eta}^{-1} \mathbf{V}_k)$, and
 $q(\mathbf{V}_k) = \mathcal{IW}(\mathbf{V}_k | \tilde{\gamma}, \tilde{\gamma} \tilde{\Sigma})$, which gives the updating algorithm for the corresponding tilded parameters.
- ▶ **\mathcal{F} evaluation:** After each E step and M step, we can also evaluate the expression of $\mathcal{F}(q)$ which can be used for **stopping rule** of the iterative algorithm.
- ▶ Final value of $\mathcal{F}(q)$ for each value of K , noted \mathcal{F}_k , can be used as a criterion for **model selection**, i.e.; **the determination of the number of clusters**.

VBA: choosing the good families for q

- ▶ Main question: We approximate $p(X)$ by $q(X)$. What are the quantities we have conserved?
 - ▶ a) Modes values: $\arg \max_x \{p(X)\} = \arg \max_x \{q(X)\}$?
 - ▶ b) Expected values: $E_p(X) = E_q(X)$?
 - ▶ c) Variances: $V_p(X) = V_q(X)$?
 - ▶ d) Entropies: $H_p(X) = H_q(X)$?
- ▶ Recent works shows some of these under some conditions.
- ▶ For example, if $p(x) = \frac{1}{Z} \exp [-\phi(x)]$ with $\phi(x)$ convex and symmetric, properties a) and b) are satisfied.
- ▶ Unfortunately, this is not the case for variances or other moments.
- ▶ If p is in the exponential family, then choosing appropriate conjugate priors, the structure of q will be the same and we can obtain appropriate **fast optimization algorithms**.

Conclusions

- ▶ Bayesian approach with Hierarchical prior model with hidden variables are very powerful tools for inverse problems and Machine Learning.
- ▶ The computational cost of all the sampling methods (MCMC and many others) are too high to be used in practical high dimensional applications.
- ▶ We explored VBA tools for effective approximate Bayesian computation.
- ▶ Application in different inverse problems in imaging system (3D X ray CT, Microwaves, PET, Ultrasound, Optical Diffusion Tomography (ODT), Acoustic source localization,...)
- ▶ Clustering and classification of a set of data are between the most important tasks in statistical researches for many applications such as data mining in biology.
- ▶ Mixture models are classical models for these tasks.
- ▶ We proposed to use a mixture of generalised Student-t distribution model for more robustness.
- ▶ To obtain fast algorithms and be able to handle large data