

Bayesian inference for inverse problems

Ali Mohammad-Djafari

*Laboratoire des Signaux et Systèmes,
Supélec, Plateau de Moulon, 91192 Gif-sur-Yvette, France*

Abstract. Traditionally, the MaxEnt workshops start by a tutorial day. This paper summarizes my talk during 2001's workshop at John Hopkins University. The main idea in this talk is to show how the Bayesian inference can naturally give us all the necessary tools we need to solve real inverse problems: starting by simple inversion where we assume to know exactly the forward model and all the input model parameters up to more realistic advanced problems of myopic or blind inversion where we may be uncertain about the forward model and we may have noisy data.

Starting by an introduction to inverse problems through a few examples and explaining their ill posedness nature, I briefly presented the main classical deterministic methods such as data matching and classical regularization methods to show their limitations. I then presented the main classical probabilistic methods based on likelihood, information theory and maximum entropy and the Bayesian inference framework for such problems. I show that the Bayesian framework, not only generalizes all these methods, but also gives us natural tools, for example, for inferring the uncertainty of the computed solutions, for the estimation of the hyperparameters or for handling myopic or blind inversion problems. Finally, through a deconvolution problem example, I presented a few state of the art methods based on Bayesian inference particularly designed for some of the mass spectrometry data processing problems.

INTRODUCTION

Forward and inverse problems

In experimental science, it is hard to find an example where we can measure directly a desired quantity. Describing *mathematical models* to relate the measured quantities to the unknown quantity of interest is called *forward modeling problem*. The main object of a forward modeling is to be able to generate data which are as likely as possible to the observed data if the unknown quantity was known. But, almost always, we want to use this model and the observed data to make inference on the unknown quantity of interest: This is the *inversion problem*. To be more explicit, let take an example that we will use all along this paper to illustrate the different aspects of inverse problems. The example is taken from the mass spectrometry where the ideal physical quantity of interest is the components mass distribution of the material under the test. There are many techniques used in mass spectrometry. The Time-of-Flight (TOF) technique is one of them. In this technique, one measures the electrical current generated on the surface of a detector by the charged ions generated by the material under the test. Finding a very fine physical model to relate the time variation of this current to the distribution of the arrival times of the charged ions, which is itself related to the components mass distribution of the material under the test, is not an easy task. However, in a first approximation, assuming

that the instrument is linear and its characteristics do not change during the acquisition time of the experiment, a very simple convolution model relates the raw data $g(t)$ to the unknown quantity of interest $f(t)$:

$$g(\tau) = \int f(t) h(\tau - t) dt, \quad (1)$$

where $h(t)$ is the point spread function (psf) of the instrument. Figure 1 shows an example of data observed (signal in b) for a theoretical mass distribution (signal in a).

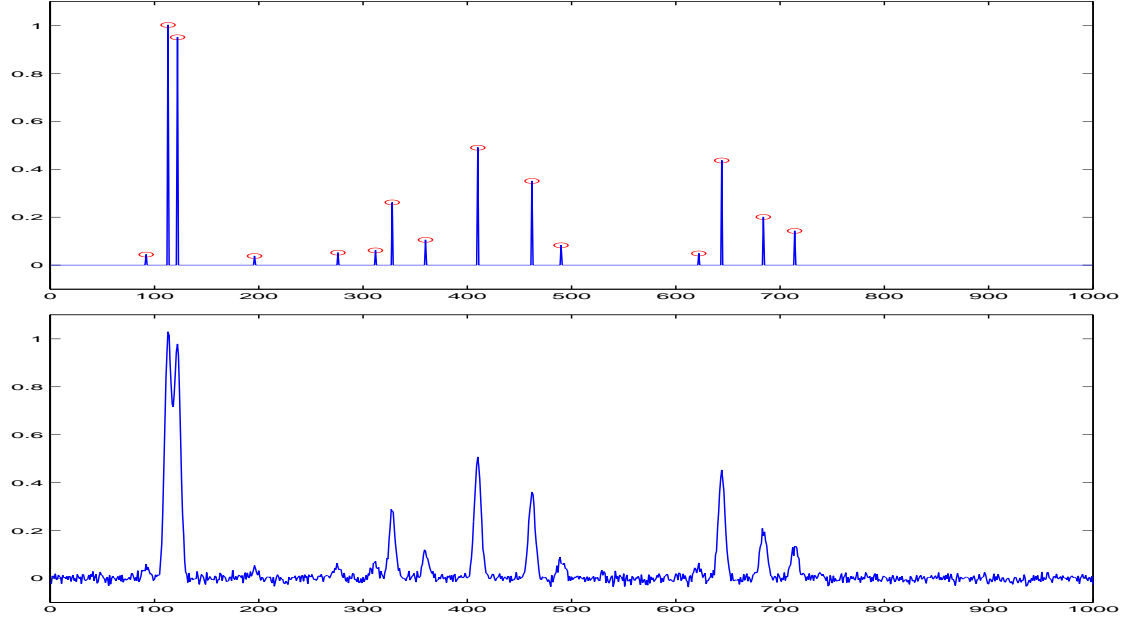


FIGURE 1. Blurring effect in TOF mass spectrometry data: a) desired or theoretical spectrum, b) observed data.

In this example, the forward problem consists in computing g given f and h which is given by a simple *convolution* operation. The inverse problem of inferring f given g and h is called *deconvolution*, the inverse problem of inferring h given g and f is called *psf identification* and the inverse problem of inferring h and f given only g is called *blind deconvolution*.

In my talk, I have given many more examples such as image restoration

$$g(x', y') = \int \int f(x, y) h(x' - x, y' - y) dx dy, \quad (2)$$

or Fourier synthesis inversion

$$g(\tau) = \int f(\omega) \exp \{-j\omega\tau\} d\omega \quad (3)$$

as well as a few non linear inverse problems. I am not going to detail them here, but I try to give a unified method to deal with all these problems. For this purpose, first we note that, in all these problems, we have always limited the number of data, for example

$y_i = g(\tau_i)$, $i = 1, \dots, m$. We also note that, to be able to do numerical computation, we need to model the unknown function f by a finite number of parameters $\mathbf{x} = [x_1, \dots, x_n]$. As an example, we may assume that

$$f(t) = \sum_{j=1}^n x_j b_j(t) \quad (4)$$

where $b_j(t)$ are known basis functions. With this assumption the raw data $\mathbf{y} = [y_1, \dots, y_m]$ are related to the unknown parameters \mathbf{x} by

$$y_i = g(\tau_i) = \sum_{j=1}^n H_{i,j} x_j \quad \text{with} \quad H_{i,j} = \int b_j(t) h(t - \tau_i) dt \quad (5)$$

which can be written in the simple matrix form $\mathbf{y} = \mathbf{H}\mathbf{x}$. The inversion problem can then be simplified to the estimation of \mathbf{x} given \mathbf{H} and \mathbf{y} . Two approaches are then in competition:

- i) the dimensional control approach which consists in an appropriate choice of the basis functions $b_j(\mathbf{r})$ and $n \leq m$ in such a way that the equation $\mathbf{y} = \mathbf{A}\mathbf{x}$ be well conditioned;
- ii) the more general regularization approach where a classical sampling basis for $b_j(\mathbf{r})$ with desired resolution is chosen no matter if $n > m$ or if \mathbf{A} is ill conditioned. In the following, we follow the second approach which is more flexible for adding more general prior information on \mathbf{x} .

We must also remark that, in general, it is very difficult to give a very fine mathematical model to take account for all the different quantities affecting the measurement process. However, we can almost always come up with a more general relation such as

$$y_i = \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}) + \epsilon_i, \quad i = 1, \dots, m \quad (6)$$

where $\boldsymbol{\theta}$ represents the unknown parameters of the forward model (for example the amplitude and the width of a Gaussian shape psf in a deconvolution problem) and $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_m]$ represents all the errors (measurement noise, discretization errors and all the other uncertainties of the model). For the case of linear models we have

$$\mathbf{y} = \mathbf{H}_{\boldsymbol{\theta}}\mathbf{x} + \boldsymbol{\epsilon}. \quad (7)$$

In this paper we focus on this general problem. We first consider the case where the model is assumed to be perfectly known. This is the simple *inversion problem*. Then we consider the more general case where we have also to infer on $\boldsymbol{\theta}$. This is the *myopic* or *blind inversion* problem.

Even in the simplest case of perfectly known linear system and exact data:

- i) the operator \mathbf{H} may not be invertible (\mathbf{H}^{-1} does not exist);
- ii) it may admit more than one inverse ($\exists \mathbf{G}_1$ and $\mathbf{G}_2 | \mathbf{G}_1(\mathbf{H}) = \mathbf{G}_2(\mathbf{H}) = \mathbf{I}$ where \mathbf{I} is the identity operator); or
- iii) it may be very ill-posed or ill-conditioned (meaning that there exists \mathbf{x} and $\mathbf{x} + \alpha\delta\mathbf{x}$ for which $\|\mathbf{H}^{-1}(\mathbf{x}) - \mathbf{H}^{-1}(\mathbf{x} + \alpha\delta\mathbf{x})\|$ never vanishes even if $\alpha \mapsto 0$ [?, ?]).

These are the three necessary conditions of *existence*, *uniqueness* and *stability* of Hadamard for the well-posedness of an inversion problem. This explains the fact that,

in general, even in this simple case, many naïve methods based on generalized inversion or on least squares may not give satisfactory results. The following figure shows, in a simple way, the ill-posedness of a deconvolution problem. On this figure, we see that three different input signals can result three outputs which are practically indistinguishable from each other. This means that, data matching alone can not distinguish between any of these inputs.

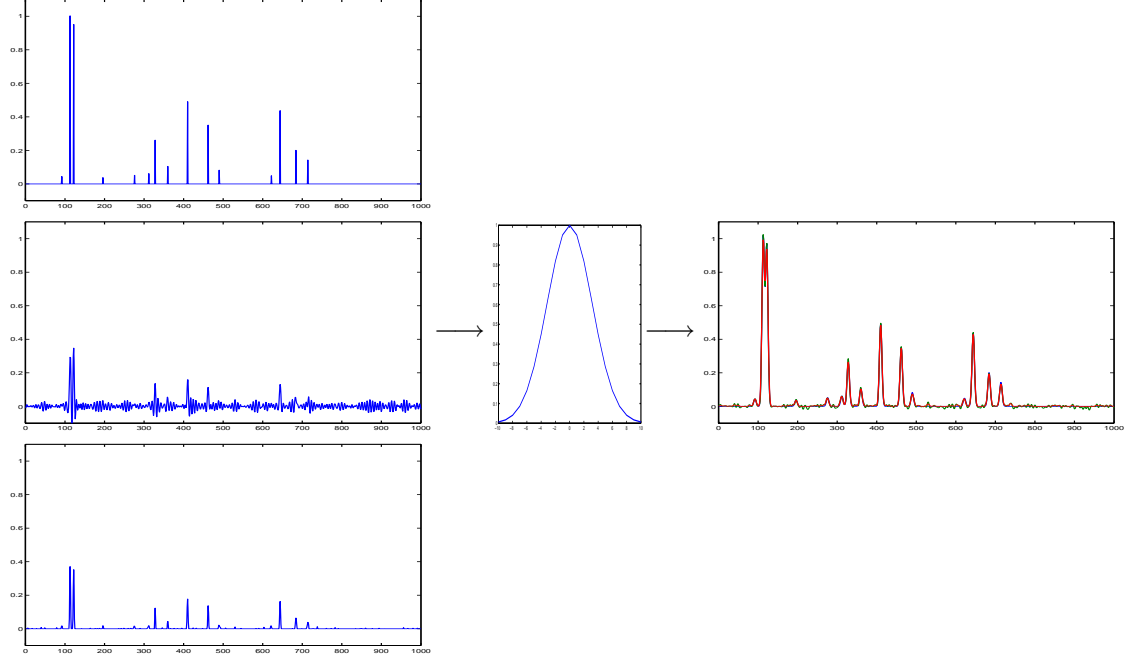


FIGURE 2. Ill-posedness of a deconvolution problem: Inputs on the left give practically indistinguishable outputs.

As a conclusion, we see that, apart from the data, we need extra information. The art of *inversion* in a particular inverse problem is how to include *just enough prior information* to obtain a satisfactory result. In the following, first we summarize the classical deterministic approaches of data matching and regularization. Then, we focus on probabilistic approaches where errors and uncertainties are taken into account through the probability laws. Here, we distinguish, three classes of methods: those which only account for the data errors (error probability distribution matching and likelihood based methods), those which only account for uncertainties of unknown parameters (entropy based methods) and those which account for both of them (Bayesian inference approach).

DATA MATCHING AND REGULARIZATION METHODS

Exact data matching

Let consider the discretized equation $y_i = h_i(\mathbf{x}) + \epsilon_i$, $i = 1, \dots, m$; and assume first that the model and data are exact ($\epsilon_i = 0$). We can then write $\mathbf{y} = \mathbf{h}(\mathbf{x})$.

Assume now the system of equations is under determined, *i.e.*, there is more than one solution satisfying it (for example when the number of data is less than the number of unknowns). Then, one way to obtain a unique solution is to define an *a priori* criterion, for example $\Delta(\mathbf{x}, \mathbf{m})$ to choose that unique solution by

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{h}(\mathbf{x})=\mathbf{y}} \{\Delta(\mathbf{x}, \mathbf{m})\} \quad (8)$$

where \mathbf{m} is an *a priori* solution and Δ a distance measure.

In the linear inverse problems case, the solution to this constrained optimization can be obtained via Lagrangian techniques which consists in defining the Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \Delta(\mathbf{x}, \mathbf{m}) + \boldsymbol{\lambda}^t(\mathbf{y} - \mathbf{H}\mathbf{x})$ and searching for $(\hat{\boldsymbol{\lambda}}, \hat{\mathbf{x}})$ through

$$\begin{cases} \hat{\boldsymbol{\lambda}} &= \arg \min_{\boldsymbol{\lambda}} \{\mathcal{D}(\boldsymbol{\lambda}) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})\} \\ \hat{\mathbf{x}} &= \arg \min_{\mathbf{x}} \{\mathcal{L}(\mathbf{x}, \hat{\boldsymbol{\lambda}})\} \end{cases} \quad (9)$$

Noting that $\nabla_{\mathbf{x}} \mathcal{L} = \nabla_{\mathbf{x}} \Delta(\mathbf{x}, \mathbf{m}) - \mathbf{H}^t \boldsymbol{\lambda}$ and $\nabla_{\boldsymbol{\lambda}} \mathcal{L} = \mathbf{y} - \mathbf{H}\mathbf{x}$ and defining $\mathcal{G}(\mathbf{s}, \mathbf{m}) = \sup_{\mathbf{x}} \{\mathbf{x}^t \mathbf{s} - \Delta(\mathbf{x}, \mathbf{m})\}$ the algorithm to find the solution $\hat{\mathbf{x}}$ becomes:

- Determine $\mathcal{G}(\mathbf{s}, \mathbf{m}) = \sup_{\mathbf{x}} \{\mathbf{x}^t \mathbf{s} - \Delta(\mathbf{x}, \mathbf{m})\}$;
- Find $\hat{\boldsymbol{\lambda}} = \arg \min_{\boldsymbol{\lambda}} \{\mathcal{D}(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^t \mathbf{y} - \mathcal{G}(\mathbf{H}^t \boldsymbol{\lambda}, \mathbf{m})\}$;
- Determine $\hat{\mathbf{x}} = \nabla_{\mathbf{s}} \mathcal{G}(\mathbf{H}^t \hat{\boldsymbol{\lambda}})$.

As an example, when $\Delta(\mathbf{x}, \mathbf{m}) = \frac{1}{2} \|\mathbf{x} - \mathbf{m}\|^2$ then $\mathcal{G}(\mathbf{s}, \mathbf{m}) = \mathbf{m}^t \mathbf{s} + \frac{1}{2} \|\mathbf{s}\|^2$, $\nabla_{\mathbf{s}} \mathcal{G} = \mathbf{m} + \mathbf{s}$ and $\mathcal{D}(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^t \mathbf{y} - \mathbf{m}^t \mathbf{H}^t \boldsymbol{\lambda} + \frac{1}{2} \|\mathbf{H}^t \boldsymbol{\lambda}\|^2$ which results to $\hat{\boldsymbol{\lambda}} = (\mathbf{H}\mathbf{H}^t)^{-1}(\mathbf{y} - \mathbf{H}\mathbf{m})$ and the solution is given by

$$\hat{\mathbf{x}} = \mathbf{m} + \mathbf{H}^t (\mathbf{H}\mathbf{H}^t)^{-1} (\mathbf{y} - \mathbf{H}\mathbf{m}). \quad (10)$$

One can remark that, when $\mathbf{m} = \mathbf{0}$ we have $\hat{\mathbf{x}} = \mathbf{H}^t (\mathbf{H}\mathbf{H}^t)^{-1} \mathbf{y}$ and this is the classical minimum norm generalized inverse solution.

Another example is the classical *Maximum Entropy* method case where $\Delta(\mathbf{x}, \mathbf{m}) = \text{KL}(\mathbf{x}, \mathbf{m})$ is the Kullback-Leibler distance or cross entropy between \mathbf{x} and the *a priori* solution \mathbf{m} :

$$\text{KL}(\mathbf{x}, \mathbf{m}) = \sum_j x_j \ln \frac{x_j}{m_j} - (x_j - m_j) \quad (11)$$

Here, the solution is given by

$$\hat{x}_j = m_j \exp \left[-[A^t \hat{\boldsymbol{\lambda}}]_j \right] \text{ with } \hat{\boldsymbol{\lambda}} = \arg \min_{\boldsymbol{\lambda}} \{\mathcal{D}(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^t \mathbf{y} - \mathcal{G}(A^t \boldsymbol{\lambda}, \mathbf{m})\} \quad (12)$$

where $\mathcal{G}(\mathbf{s}, \mathbf{m}) = \sum_j m_j (1 - \exp[-s_j])$. But, unfortunately here $\mathcal{D}(\boldsymbol{\lambda})$ is not a quadratic function of $\boldsymbol{\lambda}$ and thus there is not an analytic expression for $\hat{\boldsymbol{\lambda}}$. However, it can be computed numerically and many algorithms have been proposed for its efficient computation. See for example [?] and the cited references for more discussions on the computational issues and algorithm implementation.

The main issue here is that, this approach gives a satisfactory solution to the uniqueness of the inverse problem, but in general, the performances obtained by the resulting algorithms stay sensitive to error on the data.

Least squares data matching and regularization

When the discretized equation $\mathbf{y} = \mathbf{h}(\mathbf{x})$ is over-determined, *i.e.*, there is no solution satisfying it exactly (for example when the number of data is greater than the number of unknowns or when the data are not exact), one can try to estimate them by:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ \Delta(\mathbf{y}, \mathbf{h}(\mathbf{x})) \}, \quad (13)$$

where $\Delta(\mathbf{y}, \mathbf{h}(\mathbf{x}))$ is a distance measure in the data space. The case where $\Delta(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \|\mathbf{y} - \mathbf{h}(\mathbf{x})\|^2$ is the classical Least Squares (LS) criterion.

For a linear inversion problem $\mathbf{y} = \mathbf{H}\mathbf{x}$, it is easy to see that any $\hat{\mathbf{x}}$ which satisfies the normal equation $\mathbf{H}^t \mathbf{H} \hat{\mathbf{x}} = \mathbf{H}^t \mathbf{y}$ is a LS solution. If $\mathbf{H}^t \mathbf{H}$ is invertible and well-conditioned then $\hat{\mathbf{x}} = (\mathbf{H}^t \mathbf{H})^{-1} \mathbf{H}^t \mathbf{y}$ is again the unique generalized inverse solution. But, in general, this is not the case: $\mathbf{H}^t \mathbf{H}$ is rank deficient and we need to constrain the space of the admissible solutions. The constraint LS is then defined as

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{C}} \{ \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 \}. \quad (14)$$

where \mathcal{C} is a convex set. The choice of the set \mathcal{C} is primordial to satisfy the three conditions of a well-posed solution. An example is the positivity constraint: $\mathcal{C} = \{\mathbf{x} : \forall j, x_j > 0\}$. Another example is $\mathcal{C} = \{\mathbf{x} : \|\mathbf{x}\| < \alpha\}$ where the solution can be computed via the optimization of

$$J(\mathbf{x}) = \|\mathbf{y} - \mathbf{H}(\mathbf{x})\|^2 + \lambda \|\mathbf{x}\|. \quad (15)$$

The main technical difficulty is the relation between α and λ . The minimum norm LS solution can also be computed using the singular value decomposition [?]. The main issue here is that, even if this approach has been well understood and commonly used, it assumes implicitly that the noise and the \mathbf{x} are Gaussian. This may not be suitable in some applications, and more specifically in mass spectrometry data processing where the unknowns are spiky spectra.

A more general regularization procedure is to define the solution to the inversion problem $\mathbf{y} = \mathbf{H}(\mathbf{x}) + \epsilon$ as the optimizer of a compound criterion $J(\mathbf{x}) = \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \lambda \phi(\mathbf{x})$ or the more general criterion

$$J(\mathbf{x}) = \Delta_1(\mathbf{y}, \mathbf{H}\mathbf{x}) + \lambda \Delta_2(\mathbf{x}, \mathbf{m}). \quad (16)$$

where Δ_1 and Δ_2 are two distances or discrepancy measures, λ a regularization parameter and \mathbf{m} an *a priori* solution[?]. The main questions here are: i) how to choose Δ_1 and Δ_2 and ii) how to determine λ and \mathbf{m} .

For the first question, many choices exist:

- Quadratic or L_2 distance: $\Delta(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|^2 = \sum_j (x_j - z_j)^2$;
- L_p distance: $\Delta(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|^p = \sum_j |x_j - z_j|^p$;
- Kullback distance: $\Delta(\mathbf{x}, \mathbf{z}) = \sum_j x_j \ln(x_j/z_j) - (x_j - z_j)$;
- Roughness distance: $\Delta(\mathbf{x}, \mathbf{z})$ any of the previous distances with $z_j = x_{j-1}$ or $z_j = (x_{j-1} + x_{j+1})/2$ or any linear function $z_j = \psi(x_k, k \in \mathcal{N}(j))$ where $\mathcal{N}(j)$ stands for the neighborhood of j . (One can see the link between this last case and the Gibbsian energies in the Markovian modeling of signals and images.)

The second difficulty in this deterministic approach is the determination of the regularization parameter λ . Even if there are some techniques based on cross validation [? , ? , ?], there is not natural tools for their extension to other hyperparameters in a natural way.

As a simple example, we consider the case where both Δ_1 and Δ_2 are quadratic: $J(\mathbf{x}) = \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\mathbf{W}}^2 + \lambda \|\mathbf{x} - \mathbf{m}\|_{\mathbf{Q}}^2$. The optimization problem, in this case, has an analytic solution:

$$\hat{\mathbf{x}} = (\mathbf{H}^t \mathbf{W} \mathbf{H} + \lambda \mathbf{Q})^{-1} (\mathbf{H}^t \mathbf{W} \mathbf{y} - \mathbf{Q} \mathbf{m}) \quad (17)$$

which can also be written

$$\hat{\mathbf{x}} = \mathbf{m} + \mathbf{Q}^{-1} \mathbf{H}^t (\mathbf{H} \mathbf{Q}^{-1} \mathbf{H}^t + \lambda^{-1} \mathbf{W}^{-1})^{-1} (\mathbf{y} - \mathbf{H} \mathbf{m}) \quad (18)$$

which is a linear function of the *a priori* solution \mathbf{m} and the data \mathbf{y} . Note also that when $\mathbf{m} = \mathbf{0}$, $\mathbf{Q} = \mathbf{I}$ and $\mathbf{W} = \mathbf{I}$ we have $\hat{\mathbf{x}} = (\mathbf{H}^t \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^t \mathbf{y}$ or $\hat{\mathbf{x}} = \mathbf{H}^t (\mathbf{H} \mathbf{H}^t + \lambda^{-1} \mathbf{I})^{-1} \mathbf{y}$ and when $\lambda = 0$ we obtain the generalized inverse solutions $\hat{\mathbf{x}} = (\mathbf{H}^t \mathbf{H})^{-1} \mathbf{H}^t \mathbf{y}$ or $\hat{\mathbf{x}} = \mathbf{H}^t (\mathbf{H} \mathbf{H}^t)^{-1} \mathbf{y}$.

As we mentioned before, the main practical difficulties in this approach are the choice of Δ_1 and Δ_2 and the determination of the hyperparameters λ and the inverse covariance matrices \mathbf{W} and \mathbf{Q} .

As a main conclusion on these deterministic inversion methods, we can say that, even if, in practice, they are used and give satisfaction, they lack tools to handle with uncertainties and to account for more precise *a priori* knowledge of statistical properties of errors and unknown parameters. The probabilistic methods can exactly handle more easily these problems as we will see in the following.

PROBABILISTIC METHODS

Probability distribution matching and maximum likelihood

The main idea here is to account for data and model uncertainty through the assignment of a theoretical distribution $p_{Y|X}(\mathbf{y}|\mathbf{x})$ to the data. In probability distribution matching method, the main idea is to determine the unknown parameters \mathbf{x} by minimizing a distance measure $\Delta(\rho, p)$ between the empirical histogram ρ of the data defined

as

$$\rho(\mathbf{z}) \stackrel{\Delta}{=} \frac{1}{N} \sum_i \delta(z_i - y_i) \quad (19)$$

and the theoretical distribution of the data $p_{Y|X}(\mathbf{z}|\mathbf{x})$.

When $\Delta(p, \rho)$ is chosen to be the Kullback-Leibler mismatch measure

$$\begin{aligned} KL[\rho, p] &\stackrel{\Delta}{=} \int \rho(\mathbf{z}) \ln \frac{\rho(\mathbf{z})}{p_{Y|X}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= - \int \rho(\mathbf{z}) \ln p_{Y|X}(\mathbf{z}|\mathbf{x}) d\mathbf{z} + \int \rho(\mathbf{z}) \ln \rho(\mathbf{z}) d\mathbf{z} \end{aligned} \quad (20)$$

we have

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{KL[\rho, p]\} = \arg \min_{\mathbf{x}} \left\{ - \int \rho(\mathbf{z}) \ln p_{Y|X}(\mathbf{z}|\mathbf{x}) d\mathbf{z} \right\}. \quad (21)$$

It is then easy to see that, for the i.i.d. data, this estimate becomes equivalent to the maximum likelihood (ML) estimate

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \{ - \ln p_{Y|X}(\mathbf{z}|\mathbf{x}) |_{\mathbf{z}=\mathbf{y}} \} = \arg \max_{\mathbf{x}} \{ p_{Y|X}(\mathbf{y}|\mathbf{x}) \}. \quad (22)$$

In the case of a linear model and Gaussian noise, it is easy to show that the ML estimate becomes equivalent to the LS one, which in general, does not give satisfactory results as we have discussed it in the previous section.

The important point to note here is that, in this approach, only the data uncertainty is considered and modeled through the probability law $p_{Y|X}(\mathbf{y}|\mathbf{x})$. We will see in the following that, in contrary to this approach, in information theory and maximum entropy methods, the data and model are assumed to be exact and only the uncertainty of \mathbf{x} is modeled through an *a priori* reference measure $\mu(\mathbf{x})$ which is updated to an *a posteriori* probability law $p(\mathbf{x})$ by optimizing the KL mismatch $KL(p, \mu)$ subject to the data constraints.

Maximum entropy in the mean

The main idea in this approach is to consider \mathbf{x} as the mean value of a quantity $\mathbf{X} \in \mathcal{C}$, where \mathcal{C} is a compact set on which we want to define a probability law $P: \mathbf{x} = \mathbb{E}_P \{ \mathbf{X} \}$ and the data \mathbf{y} as exact equality constraints on it:

$$\mathbf{y} = \mathbf{H}\mathbf{x} = \mathbf{H}\mathbb{E}_P \{ \mathbf{X} \} = \int_{\mathcal{C}} \mathbf{H}\mathbf{x} dP(\mathbf{x}). \quad (23)$$

Then, assuming that we can translate our prior information on the unknowns through a prior law (a reference measure) $d\mu(\mathbf{x})$, we can determine the distribution P by:

$$\text{maximize} \quad - \int_{\mathcal{C}} \ln \frac{dP(\mathbf{x})}{d\mu(\mathbf{x})} dP(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{y} = \mathbf{H}\mathbf{x} = \mathbf{H}\mathbb{E}_P \{ \mathbf{X} \}. \quad (24)$$

The solution is obtained via the Lagrangian:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \int_{\mathcal{C}} \left[\ln \frac{dP(\mathbf{x})}{d\mu(\mathbf{x})} - \boldsymbol{\lambda}^t (\mathbf{y} - \mathbf{H}\mathbf{x}) \right] dP(\mathbf{x})$$

and is given by: $dP(\mathbf{x}, \boldsymbol{\lambda}) = \exp [\boldsymbol{\lambda}^t [\mathbf{H}\mathbf{x}] - \ln Z(\boldsymbol{\lambda})] d\mu(\mathbf{x})$, where

$Z(\boldsymbol{\lambda}) = \int_{\mathcal{C}} \exp [\boldsymbol{\lambda}^t [\mathbf{H}\mathbf{x}]] d\mu(\mathbf{x})$. The Lagrange parameters are obtained by searching the unique solution (if exists) of the following system of non linear equations:

$$\frac{\partial \ln Z(\boldsymbol{\lambda})}{\partial \lambda_i} = y_i, \quad i = 1, \dots, M. \quad (25)$$

Then, naturally, the solution to the inverse problem is defined as the expected value of this distribution: $\hat{\mathbf{x}}(\boldsymbol{\lambda}) = E_P \{ \mathbf{X} \} = \int \mathbf{x} dP(\mathbf{x}, \boldsymbol{\lambda})$. The interesting point here is that, the solution $\hat{\mathbf{x}}(\hat{\boldsymbol{\lambda}})$ can be computed without actually computing P in two ways:

– Via optimization of a dual criterion: The solution $\hat{\mathbf{x}}$ is expressed as a function of the dual variable $\hat{\mathbf{s}} = \mathbf{H}^t \hat{\boldsymbol{\lambda}}$ by $\hat{\mathbf{x}}(\hat{\mathbf{s}}) = \nabla_{\mathbf{s}} G(\hat{\mathbf{s}}, \mathbf{m})$ where

$$G(\mathbf{s}, \mathbf{m}) = \ln Z(\mathbf{s}, \mathbf{m}) = \ln \int_{\mathcal{C}} \exp [\mathbf{s}^t \mathbf{x}] d\mu(\mathbf{x}), \quad \mathbf{m} = E_{\mu} \{ \mathbf{X} \} = \int_{\mathcal{C}} \mathbf{x} d\mu(\mathbf{x})$$

and $\hat{\boldsymbol{\lambda}} = \arg \max_{\boldsymbol{\lambda}} \{ D(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^t \mathbf{y} - G(\mathbf{H}^t \boldsymbol{\lambda}) \}$.

– Via optimization of a primal or direct criterion:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{C}} \{ H(\mathbf{x}, \mathbf{m}) \} \quad \text{s.t. } \mathbf{y} = \mathbf{H}\mathbf{x} \quad \text{where } H(\mathbf{x}, \mathbf{m}) = \sup_{\mathbf{s}} \{ \mathbf{s}^t \mathbf{x} - G(\mathbf{s}, \mathbf{m}) \}.$$

Another interesting point is the link between these two options:

- i) Functions G and H depend on the reference measure $\mu(\mathbf{x})$;
- ii) The dual criterion $D(\boldsymbol{\lambda})$ depends on the data and the function G ;
- iii) The primal criterion $H(\mathbf{x}, \mathbf{m})$ is a distance measure between \mathbf{x} and \mathbf{m} which means: $H(\mathbf{x}, \mathbf{m}) \geq 0$ and $H(\mathbf{x}, \mathbf{m}) = 0$ iff $\mathbf{x} = \mathbf{m}$; $H(\mathbf{x}, \mathbf{m})$ is differentiable and convex on \mathcal{C} and $H(\mathbf{x}, \mathbf{m}) = \infty$ if $\mathbf{x} \notin \mathcal{C}$;

iv) If the reference measure is separable: $\mu(\mathbf{x}) = \prod_{j=1}^N \mu_j(x_j)$ then P is too:

$dP(\mathbf{x}, \boldsymbol{\lambda}) = \prod_{j=1}^N dP_j(x_j, \lambda_j)$ and we have

$$G(\mathbf{s}, \mathbf{m}) = \sum_j g_j(s_j, m_j), \quad H(\mathbf{x}, \mathbf{m}) = \sum_j h_j(x_j, m_j), \quad \hat{x}_j = g'_j(s_j, m_j).$$

where g_j is the log Laplace transform (Cramer transform) of μ_j :

$$g_j(s) = \ln \int \exp [sx] d\mu_j(x);$$

and h_j is the convex conjugate of g_j : $h_j(x) = \max_{\mathbf{s}} \{ sx - g_j(s) \}$.

The following table gives three examples of choices for μ_j and the resulting expressions for g_j and h_j :

	$\mu_j(x)$	$g_j(s)$	$h_j(x, m)$
Gaussian:	$\exp [-(1/2)(x-m)^2]$	$(1/2)(s-m)^2$	$(1/2)(x-m)^2$
Poisson:	$(m^x/x!) \exp [-m]$	$\exp [m-s]$	$-x \ln (x/m) + m - x$
Gamma:	$x^{\alpha-1} \exp [-(x/m)]$	$\ln (s-m)$	$-\ln (x/m) + (x/m) - 1$

We may remark that the two famous expressions of the Burg $\ln x$ and Shannon $-x \ln x$ entropies are obtained as special cases.

As a conclusion, we see that the Maximum entropy in mean extends in some way the classical ME approach by giving other expressions for the criterion to optimize. Indeed, it can be shown that when we optimize a convex criterion subject to the data constraints we are optimizing the entropy of some quantity related to the unknowns and *vice versa*. However, as we have mentioned, basically, in this approach the data and the model are assumed to be exact even if some extensions to the approach gives the possibility to account for the errors [?]. In the next section, we see how the Bayesian approach can naturally account for both uncertainties on the data and on the unknown parameters \mathbf{x} .

BAYESIAN INFERENCE APPROACH

In Bayesian approach, the main idea is to translate our prior knowledge on the errors ϵ and on the unknowns \mathbf{x} to prior probability laws $p(\epsilon)$ and $p(\mathbf{x})$. The next step is to use the forward model and $p(\epsilon)$ to deduce $p(\mathbf{y}|\mathbf{x})$. The Bayes rule can then be used to determine the posterior law of the unknowns $p(\mathbf{x}|\mathbf{y})$ from which we can deduce any information about the unknowns \mathbf{x} . The posterior $p(\mathbf{x}|\mathbf{y})$ is thus the final product of the Bayesian approach. However, very often, we need a last step which is to take out the necessary information about \mathbf{x} from this posterior. The tools for this last step are the decision and estimation theories.

To illustrate this, let consider the case of linear inverse problems $\mathbf{y} = \mathbf{H}\mathbf{x} + \epsilon$. The first step is to write down explicitly our hypothesis: starting by the hypothesis that ϵ is zero-mean (no systematic error), white (no time correlation for the errors) and assuming that we may only have some idea about its energy $\sigma_\epsilon^2 = 1/(2\phi_1)$, and using either the intuition or the Maximum Entropy Principle (MEP) lead to a Gaussian prior law: $\epsilon \sim \mathcal{N}(\mathbf{0}, 1/(2\phi_1)\mathbf{I})$. Then, using the forward model with this assumption leads to:

$$p(\mathbf{y}|\mathbf{x}, \phi_1) \propto \exp [-\phi_1 \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2]. \quad (26)$$

The next step is to assign a prior law to the unknowns \mathbf{x} . This step is more difficult and needs more caution. In inverse problems, as we presented, \mathbf{x} represents the samples of a signal or the pixel values of an arian image. Very often then we have ensemblist prior knowledge about the signals or images concerned by the application and we can model them. The art of the engineer is then to choose the appropriate model and to translate this information to a probability law to reflect it.

Again here, let illustrate this step, first through a few general examples and then more specifically the case of mass spectrometry deconvolution problem.

In the first example, we assume that, *a priori* we do not have (or we do not want or we are not able to account for) any knowledge about the correlation between the components of \mathbf{x} . This leads us to

$$p(\mathbf{x}) = \prod_j p_j(x_j). \quad (27)$$

Now, we have to assign $p_j(x_j)$. For this, we may assume to know the mean values m_j and some idea about the dispersions about these mean values. This again leads us to Gaussian laws $\mathcal{N}(m_j, \sigma_{x_j}^2)$, and if we assume the same dispersions $\sigma_{x_j}^2 = 1/(2\phi_2), \forall j$ we obtain

$$p(\mathbf{x}) \propto \exp \left[-\phi_2 \sum_j |x_j - m_j|^2 \right] = \exp \left[-\phi_2 \|\mathbf{x} - \mathbf{m}\|^2 \right] \quad (28)$$

With these assumptions, using the Bayes rule, we obtain

$$p(\mathbf{x}|\mathbf{y}) \propto \exp \left[-\phi_1 \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 - \phi_2 \|\mathbf{x} - \mathbf{m}\|^2 \right]. \quad (29)$$

This posterior law contains all the information we can have on \mathbf{x} (combination of our prior knowledge and data). If \mathbf{x} was a scalar or a vector of only two components, we could plot the probability distribution and look at it. But, in practical applications, \mathbf{x} may be a vector with huge number of components. Then, even if we can obtain an expression for this posterior, we may need to summarize its information content. In general then, we may choose, equivalently, between summarizing it by its mode, mean, marginal modes, *etc* ..., or use the decision and estimation theory to define *point estimators* to be used to compute (*best representing values*). For example, we can choose the value $\hat{\mathbf{x}}$ which corresponds to the mode of $p(\mathbf{x}|\mathbf{y})$ – the *Maximum a posteriori (MAP)* estimate, or the value $\hat{\mathbf{x}}$ which corresponds to the mean of this posterior– the *Posterior mean (PM)* estimate, or when interested to the component x_j , to choose \hat{x}_j corresponding to the mode of the posterior marginal $p(x_j|\mathbf{y})$.

We can also generate samples from this posterior and just look at them as a movie or use them to compute the PM estimate. We can also use it to compute the posterior covariance matrix ($\mathbf{P} = \mathbb{E}\{(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^t\}$ where $\hat{\mathbf{x}}$ is the posterior mean), from which we can infer on the uncertainty of the proposed solutions.

In the Gaussian priors case we just presented, it is easy to see that, the posterior law is also Gaussian and all these estimates are the same and can be computed by minimizing

$$J(\mathbf{x}) = -\ln p(\mathbf{x}|\mathbf{y}) = \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \lambda \|\mathbf{x} - \mathbf{m}\|^2 \quad \text{with } \lambda = \frac{\phi_2}{\phi_1} = \frac{\sigma_\epsilon^2}{\sigma_x^2}. \quad (30)$$

We may note here the analogy with the quadratic regularization criterion (16) with the emphasis that the choice $\Delta_1(\mathbf{y}, \mathbf{H}\mathbf{x}) = \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2$ and $\Delta_2(\mathbf{x}, \mathbf{m}) = \|\mathbf{x} - \mathbf{m}\|^2$ are the direct consequences of Gaussian choices for prior laws of the noise $p(\epsilon)$ and the unknowns $p(\mathbf{x})$.

The Gaussian choice for $p_j(x_j)$ is not always a pertinent one. For example, we may *a priori* know that the distribution of x_j around their means m_j are more concentrated

but great deviations from them are also more likely than a Gaussian distribution. This knowledge can be translated by choosing a Generalized Gaussian law:

$$p(x_j) \propto \exp \left[-\frac{1}{2\sigma_x^2} |x_j - m_j|^p \right], \quad 1 \leq p \leq 2. \quad (31)$$

In some cases we may know more, for example we may know that x_j are positive values. Then a Gamma prior law

$$p(x_j) = \mathcal{G}(\alpha, m_j) \propto (x_j/m_j)^{-\alpha} \exp[-x_j/m_j] \quad (32)$$

would be a better choice.

In some other cases we may know that x_j are discrete positive values. Then a Poisson prior law

$$p(x_j) \propto \frac{m_j^{x_j}}{x_j!} \exp[-m_j] \quad (33)$$

is a better choice.

In all these cases, the expression of the posterior is $p(\mathbf{x}|\mathbf{y}) \propto \exp[-J(\mathbf{x})]$ with $J(\mathbf{x}) = \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \lambda\phi(\mathbf{x})$ where $\phi(\mathbf{x}) = -\ln p(\mathbf{x})$. It is interesting to note the different expressions of $\phi(\mathbf{x})$ for the prior laws discussed and remark that they contain different *entropy expressions* for the \mathbf{x} .

The last general example is the case where *a priori* we know that x_j are not independent, for example when they represents the pixels of an aerian image. We may then use a Markovian modeling

$$p(x_j|x_k, k \in \mathcal{S}) = p(x_j|x_k, k \in \mathcal{N}(j)), \quad (34)$$

where $\mathcal{S} = \{1, \dots, N\}$ stands for the whole set of pixels and $\mathcal{N}(j) = \{k : |k - j| \leq r\}$ stands for r -th order neighborhood of j .

With some assumptions on the border limits, such models again result to the optimization of the same criterion with

$$\phi(\mathbf{x}) = \Delta_2(\mathbf{x}, \mathbf{z}) = \sum_j \phi(x_j, z_j) \text{ where } z_j = \psi(x_k, k \in \mathcal{N}(j)) \quad (35)$$

with different potential functions $\phi(x_j, z_j)$.

A simple example is the case where $z_j = x_{j-1}$ and $\phi(x_j, z_j)$ any function in between the following:

$$\left\{ |x_j - z_j|^\alpha, \quad \alpha \ln \frac{x_j}{z_j} + \frac{x_j}{z_j}, \quad x_j \ln \frac{x_j}{z_j} + (x_j - z_j) \right\}$$

See ([?, ?, ?, ?]) for some more discussion and properties of these potential functions.

As one of the main conclusions here, we see that, as it concerns the MAP estimation, the Bayesian approach is equivalent to the general regularization. However, here the

choice of the distance measure $\Delta_1(\mathbf{y}, \mathbf{H}\mathbf{x})$ depends on the forward model and the hypothesis on the noise and the choice of the distance measure $\Delta_2(\mathbf{x}, \mathbf{m})$ depends on the prior law chosen for \mathbf{x} .

One more extra feature here is that, we have access to the whole posterior $p(\mathbf{x}|\mathbf{y})$ from which, not only we can define an estimate but also, we can quantify its corresponding remained uncertainty. We can also compare posterior and prior laws of the unknowns to measure the amount of information contained in the observed data. Finally, as we will see in the following, we have finer tools to model unknown signals or images and to estimate the hyperparameters.

OPEN PROBLEMS AND ADVANCED METHODS

As we have remarked in previous sections, in general, the solution of an inverse problem depends on our prior hypothesis on errors ϵ and on \mathbf{x} . Before applying the Bayes rule, we have to assign the prior laws to them. From the forward model and assumptions on ϵ we assign $p(\mathbf{y}|\mathbf{x}, \phi_1)$ and from the assumptions on \mathbf{x} we assign $p(\mathbf{x}|\phi_2)$. This step is one of the most crucial part of the applicability of the Bayesian framework for inverse problems. Modeling a signal and finding the corresponding expression for the prior law $p(\mathbf{x}|\phi_2)$ is not an easy task. This choice may have many consequences: the complexity of the computation of the posterior and consequently the computation of any point estimators such as MAP (which needs optimization) or PM (which needs integration either analytically or by Monte Carlo methods). This modeling depends also on the application. We discuss this point through the particular deconvolution problem in mass spectrometry.

Appropriate modeling of input signal

We actually had started this discussion in previous section and we saw that, at least for linear inverse problems with a white Gaussian assumption of the noise, the posterior has for expression: $p(\mathbf{x}|\mathbf{y}) \propto \exp[-J(\mathbf{x})]$ with

$$J(\mathbf{x}) = \|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 + \lambda\phi(\mathbf{x}) \quad (36)$$

with $\phi(\mathbf{x}) = -\ln p(\mathbf{x})$. Thus the expression and properties of $J(\mathbf{x})$, and consequently those of the posterior $p(\mathbf{x}|\mathbf{y})$ depend on the prior $p(\mathbf{x})$. For example if $p(\mathbf{x})$ is Gaussian then

$$\phi(\mathbf{x}) = -\ln p(\mathbf{x}) = \sum_j x_j^2$$

is a quadratic function of \mathbf{x} . Then the MAP or PM estimates have the same values and their computation needs the optimization of a quadratic criterion which can be done either analytically or by using any simple gradient based algorithm. But the Gaussian modeling is not always an appropriate one. Let take our example of deconvolution of

mass spectrometry data. We know *a priori* that the input signal must be positive. Then a truncated Gaussian will be a better choice:

$$\phi(\mathbf{x}) = \sum_j x_j^2, \quad \text{if } x_j \geq 0; \quad \text{else } \phi(\mathbf{x}) = \infty.$$

But, we know still more about the input signal: it has pulse shapes, meaning that, if we look at the histogram of the samples of a typical signal, we see that great number of samples are near to zero but great deviations from this background are not rare. Thus, a generalized Gaussian

$$\phi(\mathbf{x}) = \sum_j |x_j|^p \quad \text{with } 1 \leq p \leq 2; \quad \text{if } x_j \geq 0; \quad \text{else } \phi(\mathbf{x}) = \infty.$$

or a Gamma prior law

$$\phi(\mathbf{x}) = \sum_j \ln x_j + x_j \quad \text{if } x_j \geq 0; \quad \text{else } \phi(\mathbf{x}) = \infty.$$

would be better choices.

We can also go further in details and want to account for the fact that we are looking for atomic pulses. Then we can imagine a binary valued random vector \mathbf{z} with $p(z_j = 1) = \alpha$ and $p(z_j = 0) = 1 - \alpha$, and describe the distribution of \mathbf{x} hierarchically:

$$p(x_j | z_j) = z_j p_0(x_j) \quad (37)$$

with $p_0(x_j)$ being either a Gaussian $p(x_j) = \mathcal{N}(m, \sigma^2)$ or a Gamma law $p(x_j) = \mathcal{G}(a, b)$. The second choice is more appropriate while the first results on simpler estimation algorithms. The inference can then be done through the joint posterior

$$p(\mathbf{x}, \mathbf{z} | \mathbf{y}) \propto p(\mathbf{y} | \mathbf{x}) p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \quad (38)$$

The estimation of \mathbf{z} is then called *Detection* and that of \mathbf{x} *Estimation*. The case where we assume $p(\mathbf{z}) = \prod_j p(z_j) = \alpha^{n_1} (1 - \alpha)^{(n - n_1)}$ with n_1 the number of ones and n the length of the vector \mathbf{z} , is called Bernoulli process and this modelization for \mathbf{x} is called *Bernoulli-Gaussian* or *Bernoulli-Gamma* as a function of the choice for $p_0(x_j)$.

The difficult step in this modeling is the detection step which needs the computation of

$$p(\mathbf{z} | \mathbf{y}) \propto p(\mathbf{z}) \int p(\mathbf{y} | \mathbf{x}) p(\mathbf{x} | \mathbf{z}) d\mathbf{x} \quad (39)$$

and then its optimization over $\{0, 1\}^n$ where n is the length of the vector \mathbf{z} . The cost of the computation of the exact solution is huge (a combinatorial problem).

Many approximations to this optimization have been proposed which result to different algorithms for this detection-estimation problem [?]. Many Monte Carlo techniques have also been proposed for generating samples of \mathbf{z} and \mathbf{x} from the posterior and thus compute the PM estimates of \mathbf{x} . Giving more details on this modeling and details of corresponding algorithms is out of the scope of this paper.

The results on the following figure illustrate this discussion. Here, we used the data in figure 1 and computed \mathbf{x} by optimizing the MAP criterion (36), with different prior laws $p(\mathbf{x}) \propto \exp[-\lambda\phi(\mathbf{x})]$ in between the following choices:

- a) Gaussian: $\phi(\mathbf{x}) = \sum x_j^2$,
- b) Gaussian truncated on positive axis: $\phi(\mathbf{x}) = \sum x_j^2, x_j > 0$,
- c) Generalized Gaussian truncated on positive axis: $\phi(\mathbf{x}) = \sum |x_j|^p$ with $p = 1.1, x_j > 0$.
- d) Entropic prior $\phi(\mathbf{x}) = \sum x_j \ln x_j - x_j, x_j > 0$,
- e) Gamma prior: $\phi(\mathbf{x}) = \sum \ln x_j + x_j, x_j > 0$.

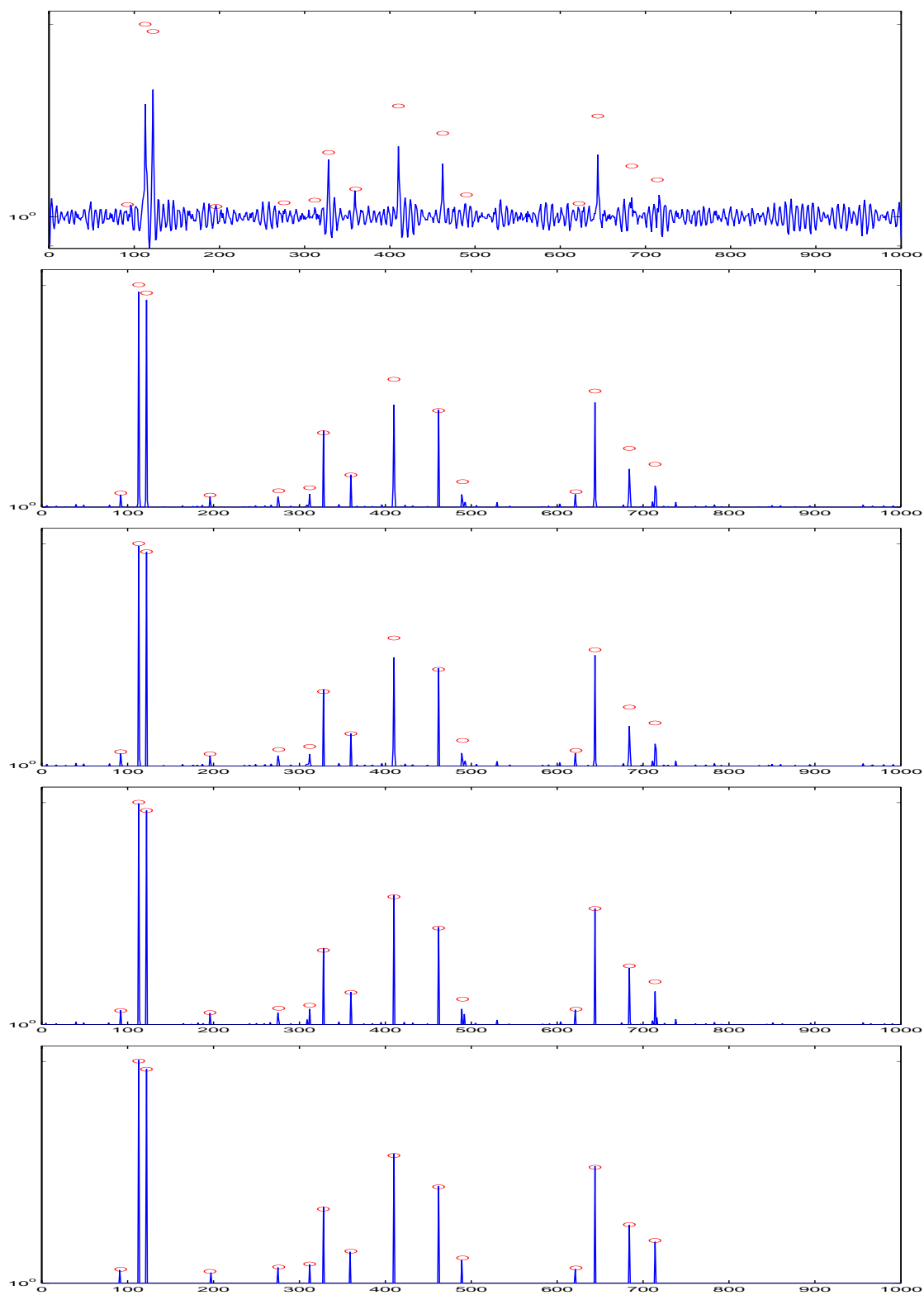


FIGURE 3. Deconvolution results with different priors: a) Gaussian b) Gaussian truncated to positive axis c) Generalized Gaussian. d) $-x \ln x$ entropic prior e) $\ln x$ entropic or Gamma prior.

As it can be seen from these results¹, for this application, the Gaussian prior does not give satisfactory result, but in almost all the other cases the results are more satisfactory, because the corresponding priors are more in agreement with the nature of the unknown input signal.

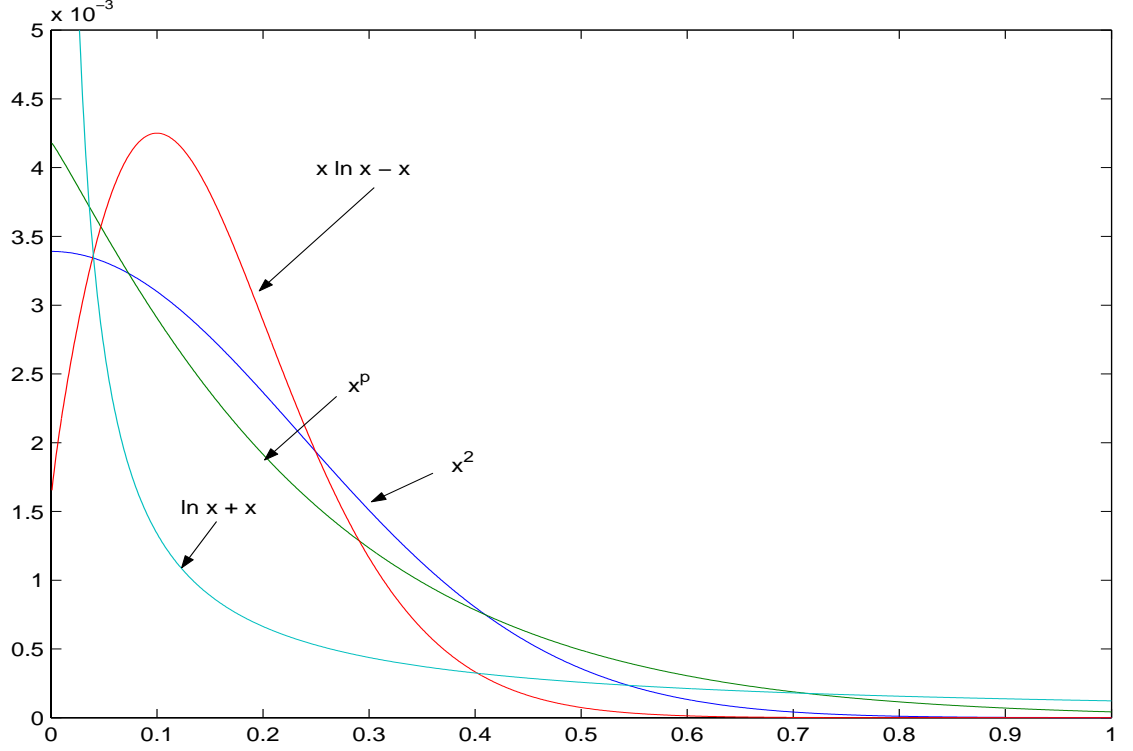


FIGURE 4. Plots of the different prior laws $p(x) \propto \exp[-\lambda\phi(x)]$: a) Truncated Gaussian $\phi(x) = x^2, \lambda = 3$ b) Truncated generalized Gaussian $\phi(x) = x^p, p = 1.1, \lambda = 4$; c) Entropic $\phi(x) = x \ln x - x, \lambda = 10$ d) Entropic $\phi(x) = \ln x + x, \lambda = 0.1$.

The Gaussian prior (a) is not at all appropriate, Gaussian truncated to positive axis (b) is a better choice. The generalized Gaussian truncated to positive axis (c) and the $-x \ln x$ entropic priors (d) give also almost the same results than the truncated Gaussian case. The Gamma prior (e) seems to give slightly better result (less missing and less artifacts) than all the others. This can be explained if we compare the shape of all these priors shown in figure (4). The Gamma prior is sharper near to zero and has longer tail than other priors. It thus favorites signals with greater number of samples near to zero and still leaves the possibility to have very high amplitude pulses. However, we must be careful on this interpretation, because all these results depend also on the hyperparameter λ whose value may be critical for this conclusion. In these experiments we used the same value for all cases. This brings us to the next open problem which is the determination of the hyperparameters.

¹ Remark that the results are presented on a logarithmic scale for the amplitudes to show in more detail the low amplitude pulses. We used $\log(1 + a)$ scale in place of y scale which has the advantage of being equal to zero for $a=0$.

Hyperparameter estimation

The Bayesian approach can be exactly applied when the direct (prior) probability laws $p(\mathbf{y}|\mathbf{x}, \phi_1)$ and $p(\mathbf{x}|\phi_2)$ are assigned. Even, when we have chosen appropriate laws, still we have to determine their parameters $\phi = [\phi_1, \phi_2]$. This problem has been addressed by many authors and the subject is an active area in statistics. See [?, ?, ?, ?], [?, ?, ?, ?] and also [?, ?, ?].

The Bayesian approach gives natural tools to handle this problem by considering $\phi = (\phi_1, \phi_2)$ as extra unknown parameters to infer on. We may then assign a prior law $p(\phi)$ to them too. However, the way to do this is also still an open problem. We do not discuss it more in this paper. The readers are invited to see [?] for some extended discussions and references. When this step is done, we can again use the Bayesian approach and compute the joint posterior $p(\mathbf{x}, \phi|\mathbf{y})$ from which we can follow three main directions:

– Joint MAP optimization: In this approach one tries to estimate both the hyperparameters and the unknown variables \mathbf{x} directly from the data by defining:

$$(\hat{\mathbf{x}}, \hat{\phi}) = \arg \max_{(\mathbf{x}, \phi)} \{p(\mathbf{x}, \phi|\mathbf{y})\} \text{ where } p(\mathbf{x}, \phi|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}, \phi) p(\mathbf{x}|\phi) p(\phi) \quad (40)$$

and where $p(\phi)$ is an appropriate prior law for ϕ . Many authors used the non informative prior law for them.

– Marginalization: The main idea in this approach is to distinguish between the two sets of unknowns: a high dimensional vector \mathbf{x} representing in general a physical quantity and a low dimensional vector ϕ representing the parameters of its prior probability laws. This argument leads to estimate first the hyperparameters by marginalizing over the unknown variables \mathbf{x} :

$$p(\phi|\mathbf{y}) \propto p(\phi) \int p(\mathbf{y}|\mathbf{x}, \phi) p(\mathbf{x}|\phi) d\mathbf{x} \quad (41)$$

and then, using them in the estimation of the unknown variables \mathbf{x} :

$$\hat{\phi} = \arg \max_{\phi} \{p(\phi|\mathbf{y})\} \longrightarrow \hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \{p(\mathbf{x}|\mathbf{y}, \hat{\phi})\}. \quad (42)$$

Note also that when $p(\phi)$ is choosed to be uniform, then $p(\phi|\mathbf{y}) \propto p(\mathbf{y}|\phi)$ which is the likelihood of the hyperparameters ϕ and the corresponding maximum likelihood (ML) estimate has all the good asymptotic properties which may not be the case for the joint MAP estimation. However, for practical applications with finite data we may not care too much about the asymptotic properties of these estimates.

– Nuisance parameters: In this approach the hyperparameters are considered as the nuisance parameters, so integrated out of $p(\mathbf{x}, \phi|\mathbf{y})$ to obtain $p(\mathbf{x}|\mathbf{y})$ and \mathbf{x} is estimated by

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \{p(\mathbf{x}|\mathbf{y})\} \text{ where } p(\mathbf{x}|\mathbf{y}) = \int p(\mathbf{y}, \mathbf{x}, \phi) d\phi \quad (43)$$

– Joint Posterior Mean: Here, \mathbf{x} and ϕ are estimated as the posterior means:

$$\hat{\mathbf{x}} = \mathbb{E}\{\mathbf{x}|\mathbf{y}\} = \int \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad \text{and} \quad \hat{\phi} = \mathbb{E}\{\phi|\mathbf{y}\} = \int \phi p(\phi|\mathbf{y}) d\phi. \quad (44)$$

The main issue here is that, excepted the first approach, all the others need integrations for which, in general, there is not analytical expressions and their numerical computation cost may be very high. At the other hand, unfortunately, the estimation by the joint maximization has not the good asymptotic properties (when number of data goes to infinity) of the estimators obtained through the marginalization or expectation. However, in finite number of data, a comparison of their relative properties is still to be done. To see some more discussions and different possible implementations of these approaches see [?]. We have also to mention that, we can always use the Markov Chain Monte Carlo (MCMC) techniques to generate samples from the joint posterior $p(\mathbf{x}, \phi|\mathbf{y})$ and then compute the joint posterior means and corresponding variances. It seems that these techniques are growing up. However, I see two main limitations for their application on real data: their huge computational cost and the need for some discussions on the tools to control their convergences.

Myopic or blind inversion problems

Consider the deconvolution problems (1) or (2) and assume now that the psf $h(t)$ or $h(x, y)$ are partially known. For example, we know they have Gaussian shape, but the amplitude a and the width σ of the Gaussian are unknown. Noting by $\theta = (a, \sigma)$ the problem then becomes the estimation of both \mathbf{x} and θ from $\mathbf{y} = \mathbf{H}_\theta \mathbf{x} + \epsilon$. The case where we know only the support of the psf but not its shape can also be casted in the same way with $\theta = [h(0), \dots, h(p)]$

Before going more in details, we must note that, in general, the blind inversion problems are much harder than the simple inversion. Taking the deconvolution problem, we have seen in introduction that, the problem even when the psf is given is ill-posed. The blind deconvolution then is still more ill-posed, because here there are more fundamental under-determinations. For example, it is easy to see that, we can find an infinite number of pairs (h, x) which result to the same convolution product $h * x$. This means that, to find satisfactory methods and algorithms for these problems need much more prior knowledge both on x and on h , and in general, the inputs must have more structures (be rich in information content) to be able to obtain satisfactory results.

Conceptually however, the problem is identical to the estimation of hyperparameters in previous section and any of the four approaches presented there can be used. One may wish however to distinguish between these parameters of the system $\theta = (a, \sigma)$ and those hyperparameters of the prior law model descriptions $\phi = (\sigma_\epsilon^2, \sigma_x^2, \dots)$. In that case, one can try to write down $p(\mathbf{x}, \theta, \phi|\mathbf{y})$ and use one of the following:

– Joint MAP estimation of \mathbf{x} , θ and ϕ : $(\hat{\mathbf{x}}, \hat{\theta}, \hat{\phi}) = \arg \max_{(\mathbf{x}, \theta, \phi)} \{p(\mathbf{x}, \theta, \phi|\mathbf{y})\}$.

– Marginalize over \mathbf{x} and estimate θ and ϕ using: $(\hat{\theta}, \hat{\phi}) = \arg \max_{(\theta, \phi)} \{p(\theta, \phi|\mathbf{y})\}$ and

then, estimate \mathbf{x} using: $\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \left\{ p(\mathbf{x}|\mathbf{y}, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\phi}}) \right\}$.

– Marginalize over \mathbf{x} and $\boldsymbol{\theta}$ and estimate $\boldsymbol{\phi}$ using: $\hat{\boldsymbol{\phi}} = \arg \max_{\boldsymbol{\phi}} \{ p(\boldsymbol{\phi}|\mathbf{y}) \}$,,

then estimate $\boldsymbol{\theta}$ using: $\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \left\{ p(\boldsymbol{\theta}|\mathbf{y}, \hat{\boldsymbol{\phi}}) \right\}$ and finally, estimate \mathbf{x} using: $\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \left\{ p(\mathbf{x}|\mathbf{y}, \hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\phi}}) \right\}$.

– Joint Posterior Mean: Here, \mathbf{x} , $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ are estimated through their respective posterior means: $\hat{\mathbf{x}} = E \{ \mathbf{x}|\mathbf{y} \}$, $\hat{\boldsymbol{\theta}} = E \{ \boldsymbol{\theta}|\mathbf{y} \}$ and $\hat{\boldsymbol{\phi}} = E \{ \boldsymbol{\phi}|\mathbf{y} \}$.

Here again, the joint optimization stays the simpler but we must be careful on interpretation of the results. For others, one can either use the Expectation-Maximization (EM) algorithms and/or MCMC sampling tools to approximately compute the necessary integration or expectation computations and overcome the computational cost issues.

CONCLUSIONS

In this paper I presented a synthetic overview of methods for inversion problems starting by deterministic data matching and regularization methods followed by a general presentation of the probabilistic methods such as error probability law matching and likelihood based and the information theory and maximum entropy based methods. Then, I focused on the Bayesian inference. I show that, as it concerns the maximum *a posteriori* estimation method, one can see easily the link with regularization methods. We discussed however the superiority of the Bayesian framework which gives naturally the necessary tools for inferring the uncertainty of the computed solution, for the estimation of the hyperparameters or for handling myopic and blind inversion problems. We saw also that probabilistic modeling of signal and images is more flexible for introduction of practical prior knowledge about them. Finally, we illustrated some of these discussions through a deconvolution example in mass spectrometry data processing.