

# JOINT ESTIMATION OF PARAMETERS AND HYPERPARAMETERS IN A BAYESIAN APPROACH OF SOLVING INVERSE PROBLEMS

Ali Mohammad-Djafari

Laboratoire des Signaux et Systèmes (CNRS-SUPELEC-UPS)  
 École Supérieure d'Électricité  
 Plateau de Moulon, 91192 Gif-sur-Yvette Cedex, France

## ABSTRACT

In this paper we propose a joint estimation of the parameters and hyperparameters (the parameters of the prior law) when a Bayesian approach with Maximum Entropy (ME) priors is used to solve the inverse problems which arise in signal and image reconstruction and restoration problems. In particular we propose two methods: one based on the Expectation Maximization (EM) algorithm who aims to find the Marginalized MAP (MMAP) estimate and the second based on a joint MAP estimation (JMAP). We discuss and compare these methods and give some simulation results in image restoration to show the relative performances of the proposed methods.

## 1. INTRODUCTION

In many signal and image reconstruction and restoration problems, the observed data  $g(\mathbf{s}_n)$  are related to the quantities of the direct physical interest  $f(\mathbf{r})$  by a linear transformation:

$$g(\mathbf{s}_m) = \int_D f(\mathbf{r})h_m(\mathbf{r})d\mathbf{r} + b(\mathbf{s}_m), \quad m = 1, \dots, M, \quad (1)$$

where  $h_m(\mathbf{r})$  is the instrument function which we assume to be known. When this integral equation is discretized, we have to solve a linear system of equations:

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \mathbf{b}, \quad (2)$$

where  $\mathbf{H}$  is a matrix whose components depend on  $h_m$ .  $\mathbf{H}$  has, in general, very large dimensions and is very ill-conditioned.  $\mathbf{f}$  is a vector whose components represents the unknown parameters to estimate (for example the pixel values of an image).  $\mathbf{g}$  is a vector whose components are the observed data, and  $\mathbf{b}$  is a vector whose components represents both the measurement noise and any other unmodeled error.

The objective of an inversion procedure is to obtain an unique and stable solution  $\hat{\mathbf{f}}$  to this *ill-posed* problem. This can be achieved by introducing prior knowledge on  $\mathbf{f}$  and on  $\mathbf{b}$  by adopting a Bayesian approach

which consists in translating our prior knowledge on  $\mathbf{f}$  and  $\mathbf{b}$  by the probability laws  $p(\mathbf{f}|\boldsymbol{\theta})$  and  $p(\mathbf{g}|\mathbf{f}, \boldsymbol{\beta})$  and using the Baye's rule to obtain  $p(\mathbf{f}|\mathbf{g}, \boldsymbol{\beta}, \boldsymbol{\theta})$ . The final step then is to choose a decision rule (for example Maximum *a posteriori* MAP) to determine a solution  $\hat{\mathbf{f}}$  to the problem.

One of the difficulties in this bayesian approach is to assign the prior law  $p(\mathbf{f}|\boldsymbol{\theta})$  in a way to reflect our prior knowledge about the solution  $\mathbf{f}$ . In cases where the prior knowledge is in the form of expectations then maximum entropy principle (MEP) can give us the solution. However, in general, no matter how  $p(\mathbf{f}|\boldsymbol{\theta})$  is assigned, it may depend on some unknown parameters  $\boldsymbol{\theta}$ , called *hyperparameters*, so that, in the inversion procedure we want to estimate also them from the data. To be more specific, let us describe a method that we developed in preceeding works [1, 2, 3] which can be resumed as follows:

- Information about noise  $\mathbf{b}$  is the covariance matrix  $\mathbf{W} = \sigma^2\mathbf{I} = (2/\beta)\mathbf{I}$ . So that, using this only information with MEP we obtain :

$$p(\mathbf{g}|\mathbf{f}, \boldsymbol{\beta}) \propto \exp[-\beta Q(\mathbf{f})] \quad \text{with} \quad Q(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 \quad (3)$$

- Information about  $\mathbf{f}$  is assumed to be in the form :

$$E\{\phi_i(\mathbf{f})\} = \mu_i \quad \text{with} \quad \phi_i(\mathbf{f}) = \sum_{n=1}^N \Omega_i(f_n), \quad i = 1, 2 \quad (4)$$

Using this information with MEP we obtain :

$$p(\mathbf{f}|\boldsymbol{\theta}) = \prod_{n=1}^N p(f_n|\boldsymbol{\theta}) \quad (5)$$

with

$$p(f_n|\boldsymbol{\theta}) = \frac{1}{Z(\theta_1, \theta_2)} \exp[-\theta_1\Omega_1(f_n) - \theta_2\Omega_2(f_n)] \quad (6)$$

and

$$Z(\theta_1, \theta_2) = \int \exp[-\theta_1\Omega_1(f) - \theta_2\Omega_2(f)] df \quad (7)$$

A scale invariance argument [3, 4] limits the set of possible functions  $\Omega_1(f)$  and  $\Omega_2(f)$  :

$$\left\{ (f^{r_1}, f^{r_2}), (f^{r_1}, \ln f), (f^{r_1}, f^{r_1} \ln f), (\ln f, \ln^2 f) \right\}. \quad (8)$$

• Applying the Bayes' rule to obtain the posterior law  $p(\mathbf{f}|\mathbf{g})$  and MAP estimation rule we obtain :

$$\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g})\} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\}, \quad (9)$$

$$\text{with } J(\mathbf{f}) = \beta Q(\mathbf{f}) + \theta_1 \phi_1(\mathbf{f}) + \theta_2 \phi_2(\mathbf{f}). \quad (10)$$

So, the well-posed Bayesian problem solving should be stated as follows:

Given  $\mathbf{H}, \mathbf{g}, \beta = 2/\sigma^2$ , and  $\{\mu_1, \mu_2\}$  or equivalently  $\{\theta_1, \theta_2\}$ , estimate  $\mathbf{f}$ .

However, in practical applications  $(\mu_1, \mu_2)$  or equivalently  $(\theta_1, \theta_2)$  are not given to us and we want also to estimate them from the data  $\mathbf{g}$ .

## 2. JOINT ESTIMATION PROBLEM

In this paper we considered the following problem :

Given  $\mathbf{H}, \mathbf{g}$  and  $\beta$  estimate  $\mathbf{f}$  and the hyperparameters  $\boldsymbol{\theta} = [\theta_1, \theta_2]$ .

Two main approaches studied and compared here are :

• **Joint Maximum a posteriori (JMAP) :**

The main idea behind this approach is to consider the hyperparameters  $\boldsymbol{\theta}$  on the same level that the other parameters  $\mathbf{f}$  and try to estimate them by

$$\text{where } (\hat{\mathbf{f}}, \hat{\boldsymbol{\theta}}) = \arg \max_{(\mathbf{f}, \boldsymbol{\theta})} \{p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g})\} \quad (11)$$

$$p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) \propto p(\mathbf{f}, \boldsymbol{\theta}, \mathbf{g}) = p(\mathbf{g}|\mathbf{f}) p(\mathbf{f}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (12)$$

Note that, we can choose either an improper or an uniform prior for  $p(\boldsymbol{\theta})$ .

This optimization problem can be implemented by successively maximizing with respect to  $\boldsymbol{\theta}$  and to  $\mathbf{f}$  :

$$\begin{cases} \hat{\boldsymbol{\theta}}^{(k)} = \arg \max_{\boldsymbol{\theta}} \{p(\hat{\mathbf{f}}^{(k)}, \boldsymbol{\theta}|\mathbf{g})\} = \arg \max_{\boldsymbol{\theta}} \{p(\hat{\mathbf{f}}^{(k)}|\boldsymbol{\theta})p(\boldsymbol{\theta})\} \\ \hat{\mathbf{f}}^{(k+1)} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}, \hat{\boldsymbol{\theta}}^{(k)}|\mathbf{g})\} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)})\} \end{cases} \quad (13)$$

Note that the first equation can be interpreted as a maximum likelihood estimate of  $\boldsymbol{\theta}$  if  $p(\boldsymbol{\theta})$  is chosen to be uniform and if  $\hat{\mathbf{f}}^{(k)}$  could be considered as a sample of the prior law  $p(\mathbf{f}|\boldsymbol{\theta})$ .

• **Marginalized MAP (MMAP) :**

The main idea behind this approach is to consider the hyperparameters  $\boldsymbol{\theta}$  on a different level than  $\mathbf{f}$ . So,  $\boldsymbol{\theta}$  is first estimated by marginalizing with respect to  $\mathbf{f}$  :

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \left\{ L(\boldsymbol{\theta}) = \int p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) d\mathbf{f} \right\} \quad (14)$$

Then using  $\hat{\boldsymbol{\theta}}$  we can estimate the solution by

$$\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}})\} \quad (15)$$

Unfortunately the analytical calculus of the integral  $L(\boldsymbol{\theta})$  is rarely (excepted for the Gaussian case) possible. However, noting that when  $p(\boldsymbol{\theta})$  is uniform  $L(\boldsymbol{\theta})$  becomes actually the likelihood and  $\hat{\boldsymbol{\theta}}$  the maximum likelihood (ML) estimate, in some cases one can obtain the solution using the Expectation-Maximisation (EM) algorithm, which in our case, is given by :

$$\begin{cases} \text{E: } Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k)}) = \mathbb{E}_{\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)}} \{\ln p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g})\} \\ \text{M: } \hat{\boldsymbol{\theta}}^{(k+1)} = \arg \max_{\boldsymbol{\theta}} \{Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k)})\} \end{cases} \quad (16)$$

## 3. COMPARISON OF JMAP AND EM-MMAP

Now, let us go a little further inside these two methods by assuming  $p(\boldsymbol{\theta})$  is chosen to be uniform. Noting back  $\boldsymbol{\theta} = (\theta_1, \theta_2)$  and replacing  $p(\mathbf{f}|\boldsymbol{\theta})$  from (4) and  $p(\mathbf{g}|\mathbf{f})$  from (3) in the equations (13) and (16) we can make the following comparison :

• **JMAP :** At iteration  $(k+1)$  of the iterative optimization algorithm (13) we have to estimate  $(\theta_1, \theta_2)^{(k+1)}$  by maximizing  $\ln p(\hat{\mathbf{f}}^{(k)}|\theta_1, \theta_2)$  with respect to  $(\theta_1, \theta_2)$ . To do this we have to deal with the following system of nonlinear equations :

$$\frac{\partial \ln Z(\theta_1, \theta_2)}{\partial \theta_i} = \frac{1}{N} \phi_i \left( \hat{\mathbf{f}}_{MAP}^{(k)} \right), \quad i = 1, 2 \quad (17)$$

where the right hand side (RHS) of these equations are

$$\frac{1}{N} \phi_i \left( \hat{\mathbf{f}}_{MAP}^{(k)} \right) = \sum_{n=1}^N \Omega_i(\hat{\mathbf{x}}_{MAP_n}^{(k)}), \quad i = 1, 2 \quad (18)$$

• **EM-MMAP :** When using the EM algorithm (16) at iteration  $(k+1)$  we have to estimate  $(\theta_1, \theta_2)^{(k+1)}$  by maximizing  $Q((\theta_1, \theta_2); (\theta_1, \theta_2)^{(k)})$  with respect to  $(\theta_1, \theta_2)$ . To do this we have to deal with the following system of nonlinear equations :

$$\frac{\partial \ln Z(\theta_1, \theta_2)}{\partial \theta_i} = \frac{1}{N} \int \phi_i(\mathbf{f}) p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)}) d\mathbf{f}, \quad i = 1, 2. \quad (19)$$

The RHS of these two equations can be written

$$RHS = \frac{1}{N} \sum_{n=1}^N \int \Omega_i(f_n) p(f_n|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)}) df_n, \quad i = 1, 2. \quad (20)$$

Thus, comparing (17) and (19), we can say that if  $p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)})$  is very concentrated around the  $\hat{\mathbf{f}}_{MAP}^{(k)}$ , i.e.,  $p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)}) \approx \delta(\mathbf{f} - \hat{\mathbf{f}}_{MAP}^{(k)})$ , then the integrals of the RHS of the equation (19) will be equivalent to RHS of (17) and the two methods will give the same numerical result. But this is not true in general.

However JMAP is easy to implement, because it does not need any integration. This is not true in the case of the EM-MMAP. Note however that, thanks to the entropic priors, the  $N$ -dimensional integration in (19) is replaced by  $N$  one-dimensional integrations in (20). However, in general, these integrals have rarely analytical solutions, so, implementing the EM algorithm is very difficult due to the need of calculating the marginal law  $p(f_n|\mathbf{g})$  and these integrals. Two solutions have been proposed in litteratures [5, 6, 7, 8] to surround this difficulty:

- 1) making a Gaussian approximation for  $p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)})$  around the maximum  $\hat{\mathbf{f}}_{MAP}$  which allows to obtain an analytic solution to these integrals, and
- 2) calculating the integrals in a stochastic method which leads us to Stochastic EM (SEM) like methods.

What we propose is an an *ad hoc* approach to approximate these integrals by

$$\int \Omega_i(f_n)p(f_n|\mathbf{g})df_n \simeq \frac{1}{k} \sum_{l=1}^k \Omega_i(f_n^{(k)}), \quad i = 1, 2. \quad (21)$$

This means that the expectations of  $\Omega_i(f_n)$  are replaced by the average values of them calculated during all the past iterations of the algorithm. Doing this, the two methods will have the same structure, i.e; to obtain new values for the  $(\theta_1, \theta_2)$  at each iteration we have to resolve the following system of equations:

$$\frac{\partial \ln Z(\theta_1, \theta_2)}{\partial \theta_i} = d_i, \quad i = 1, 2, \quad (22)$$

where

$$d_i = \begin{cases} \sum_{n=1}^N \Omega_i(\hat{x}_{MAP_n}^{(k)}) & \text{for JMAP} \\ \frac{1}{k} \sum_{l=1}^k \sum_{n=1}^N \Omega_i(\hat{x}_{MAP_n}^{(k)}) & \text{for EM-MMAP} \end{cases} \quad (23)$$

A final and very important remarque is that we have to insure that the criteria to be maximized admit at least a local maximum. This question may arises more specifically when dealing with JMAP, but it may also arises in the case of the EM-MMAP. This difficulty is indepent of the algorithm used to find the optimum. To explain more this let us have a look on the expression to optimize in JMAP (eq. 11) which can be written as

$$(\hat{\mathbf{f}}, \hat{\boldsymbol{\theta}}) = \arg \min_{(\mathbf{f}, \boldsymbol{\theta})} \{-\ln p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g})\}, \quad (24)$$

where

$$-\ln p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) = \beta \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \theta_1 \phi_1(\mathbf{f}) + \theta_2 \phi_2(\mathbf{f}) - N \ln Z(\theta_1, \theta_2) - \ln p(\boldsymbol{\theta}). \quad (25)$$

Depending on the choice of  $p(\boldsymbol{\theta})$  and  $p(\mathbf{f}|\boldsymbol{\theta})$  and consequently  $\phi_1(\mathbf{f})$ ,  $\phi_2(\mathbf{f})$  and  $Z(\theta_1, \theta_2)$ , this criterion may even not have a local minimum.

In practice this difficulty is encountered when we are dealing with the system of equations (17), (19) or (22) which are obtained by equating to zero the gradient of the above criterion with respect to  $\boldsymbol{\theta}$ . In fact this system of equations may not have any solution for given values  $d_i$ . On the same way  $\ln Z(\theta_1, \theta_2)$  which depends on the prior probability law  $p(f|\theta_1, \theta_2)$  may also be defined only on a set of admissible values of  $(\theta_1, \theta_2)$ . For example, in the Gaussian case where  $\Omega_1(f) = f^2$  and  $\Omega_2(f) = f$  we have  $Z(\theta_1, \theta_2) = \sqrt{\pi/\theta_1}$  which does not depend even on  $\theta_2$  and  $\theta_1$  must be positive.

However, when we have taken the necessary cautions, in practical situations the two algorithms will give satisfactory results. This has been shown at least in our simulations as we will see below.

#### 4. SIMULATION RESULTS

In this section, we present some simulation results which show the relative performances of the the two proposed methods in image restoration problems. For this, we have first created two synthetic (64 by 64) images (O1 and O2) (Fig. 1). According to the histograms of the images, we assigned

- a Gaussian law ( $\Omega_1(f) = f^2$ ,  $\Omega_2(f) = f$ ) to the first image, and
- a Gamma law ( $\Omega_1(f) = \ln f$ ,  $\Omega_2(f) = f$ ,  $f > 0$ ) to the second one.

Table 1. summarizes the results of parameter estimation.

images	$\Omega_1(f)$	$\Omega_2(f)$	$f$ dom.	$\theta_1$	$\theta_2$
O1	$f^2$	$f$	$\mathbf{R}$	$4.5e-4$	$-2.5e-2$
O2	$f^2$	$f$	$f > 0$	$6.9e-1$	$1.1e-2$

Table 1: Prior laws and their parameters.

In a next step, we created degraded images by blurring them with a (9 by 9) Gaussian point spread function (PSF) and added a Gaussian noise with a given variance so that the signal to noise ratio was fixed to 20dB. Fig. 2 shows the degraded images.

In the final step, using the JMAP and the EM-MMAP methods we restored these images and simultaneously estimated the hyperparameters  $(\theta_1, \theta_2)$ . Fig. 3 and Fig. 4 show the restored images by JMAP and EM-MMAP methods and the Table 2. summarizes the hyperparameter estimation results.

images	JMAP		EM-MMAP	
	$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\theta}_1$	$\hat{\theta}_2$
O1	$4.3e-4$	$-2.4e-2$	$4.3e-4$	$-2.5e-2$
O2	$6.2e-1$	$1.3e-2$	$6.3e-2$	$1.3e-2$

Table 2: Results of hyperparameter estimation and distance between the original and estimated images for the two methods JMAP and EM-MMAP.

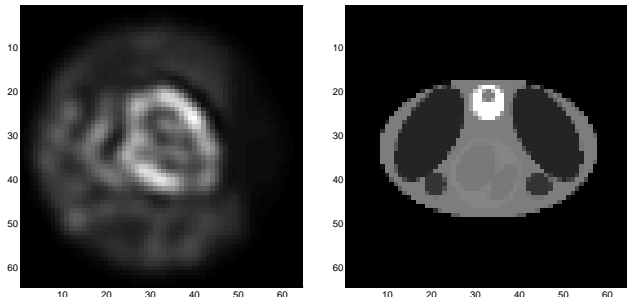


Fig. 1: Original images o1 and o2.

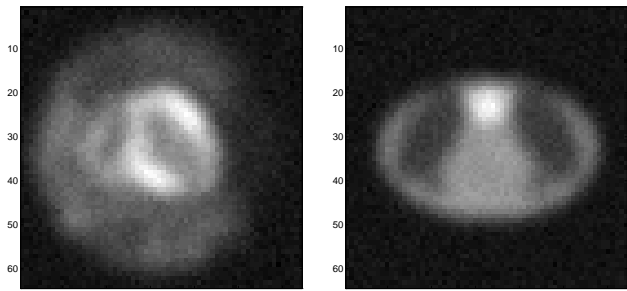


Fig. 2: Degraded images I1 and I2.

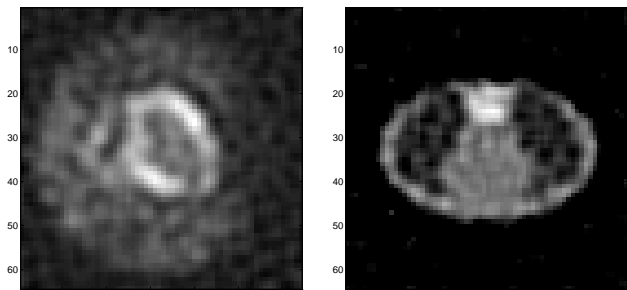


Fig. 3: Restored images by JMAP method.

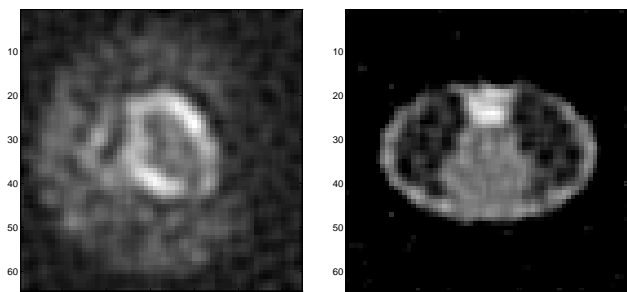


Fig. 4: Restored images by EM-MMAP method.

## 5. CONCLUSIONS

The main objective of this paper was to propose and compare two specific methods for joint estimation of

the unknowns and the hyperparameters of the inverse problems in a Bayesian approach with ME priors. The first method, named EM-MMAP, is based on the marginalised maximum likelihood (MML) and the EM algorithm and the second, named JMAP, is based on generalized maximum likelihood (GML) methods. The relative performances of these two methods have been showed in image restoration problems.

## 6. REFERENCES

- [1] A. Mohammad-Djafari and J. Idier, *Maximum Likelihood Estimation of the Lagrange Parameters of the Maximum Entropy Distributions*, pp. 131–140. Maximum entropy and Bayesian methods, Seattle, USA: Kluwer Academic Publ., smith, c.r. and erikson, g.j. and neudorfer, p.o. ed., 1991.
- [2] A. Mohammad-Djafari, “On the estimation of hyperparameters in Bayesian approach of solving inverse problems,” in *Proceedings of IEEE ICASSP*, (Minneapolis, U.S.A.), pp. 567–571, IEEE, Aprilxspace1993.
- [3] A. Mohammad-Djafari and J. Idier, “Scale invariant Bayesian estimators for linear inverse problems,” in *Proc. of the First ISBA meeting*, (San Fransisco, USA), Augustxspace1993.
- [4] S. Brette, J. Idier, and A. Mohammad-Djafari, “Scale invariant Markov models for linear inverse problems,” in *Proc. of the Section on Bayesian Statistical Sciences*, (Alicante, Spain), pp. 266–270, American Statistical Association, 1994.
- [5] A. Dempster, N. Laird, and D. Rubin, “Maximum likelihood from incomplete data via the em algorithm,” *Journal of the Royal Statistical Society B*, vol. 39, pp. 1–38, 1977.
- [6] M. Miller and D. Snyder, “The role of likelihood and entropy in incomplete-data problems: Applications to estimating point-process intensities and toeplitz constrained covariances,” *Proceedings of the IEEE*, vol. 75, pp. 892–906, Julyxspace1987.
- [7] C. F. J. Wu, “On the convergence of the em algorithm,” *Ann. Statist.*, vol. 11, no. 1, pp. 95–103, 1983.
- [8] F. Champagnat and J. Idier, “An alternative to standard maximum likelihood for Gaussian mixtures,” in *Proceedings of IEEE ICASSP*, (Detroit), pp. 2020–2023, 1995.