

Sparse Spike Train Deconvolution Using the Hunt Filter and a Thresholding Method

Vincent Mazet, David Brie, and Cyrille Caironi

Abstract—A new deconvolution method of sparse spike trains is presented. It is based on the coupling of the Hunt filter with a thresholding. We show that a good model for the probability density function of the Hunt filter output is a Gaussian mixture, from which we derive the threshold that minimizes the probability of errors. Based on an interpretation of the method as a maximum *a posteriori* (MAP) estimator, the hyperparameters are estimated using a joint MAP approach. Simulations show that this method performs well at a very low computation time.

Index Terms—Bernoulli–Gaussian, Hunt filter, sparse spike train deconvolution, thresholding.

I. INTRODUCTION

WE CONSIDER the problem of sparse spike train deconvolution which is classically stated as follows (in the sequel, v_i will represent the i th element of vector \mathbf{v} and $M_{i,j}$ the element (i, j) of the matrix \mathbf{M}): given some observation sequence $\mathbf{y} \in \mathbb{R}^N$, find the sparse spike train sequence $\mathbf{x} \in \mathbb{R}^N$ such as

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}$$

where $\mathbf{H} \in \mathbb{R}^{N \times N}$ is a known impulse response matrix, and $\mathbf{n} \in \mathbb{R}^N$ a noise term assumed to be Gaussian and independent, identically distributed (i.i.d.) whose probability density function (pdf) is $p(\mathbf{n}) = \mathcal{N}(\mathbf{0}, r_n \mathbf{I}_N)$, r_n being the noise variance. To handle the sparse nature of \mathbf{x} , it is modeled as an i.i.d. Bernoulli–Gaussian (BG) sequence [1], [2], i.e.,

$$\forall k, \quad p(\mathbf{x}_k) = \lambda \mathcal{N}(0, r_x) + (1 - \lambda) \delta(\mathbf{x}_k)$$

where r_x is the pulse variance, $\lambda \ll 1$ is the Bernoulli parameter (ratio between pulse and sample number), and δ is the Dirac mass centered on zero. In the sequel, $\boldsymbol{\theta} = \{\lambda, r_x, r_n\}$ gathers the hyperparameters of the problem.

Among the many applications of sparse spike deconvolution, we would like to mention partial discharge analysis, which has motivated this work [3]. The final goal was to develop algorithms that can be embedded to monitor the insulation of high-voltage AC motors. In addition, in this particular application,

Manuscript received April 17, 2003; revised July 3, 2003. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Dimitris A. Pados.

V. Mazet and D. Brie are with the Centre de Recherche en Automatique de Nancy, Centre National de la Recherche Scientifique, UMR 7039, Université Henri Poincaré, 54506 Vandoeuvre-lès-Nancy Cedex, France (e-mail: vincent.mazet@cran.uhp-nancy.fr; david.brie@cran.uhp-nancy.fr).

C. Caironi is with Alstom Moteurs, 54250 Champigneulle, France (e-mail: cyrille.caironi@tde.alstom.com).

Digital Object Identifier 10.1109/LSP.2004.826655

one has to process very large signals (typically $N = 2^{16}$ points). These are the reasons that have led us to develop an approach fast and simple to implement needing only low memory space.

On the one hand, combinatory algorithms like single most likely replacement (SMLR) [1], [2], iterated window maximization (IWM) [4], or others given in [5] are well-suited to the problem. Although they give very satisfactory estimations, they are very slow and difficult to implement on embedded systems.

On the other hand, the Hunt filter [6] is very fast and simple to implement, but it does not yield satisfactory results, mainly because of the underlying prior, which assumes that the restored signal is Gaussian. To overcome this drawback, we propose to associate the Hunt filter with a thresholding procedure to actually restore a sparse spike train signal resulting in the so-called Hunt–threshold (HT) method.

The letter is organized as follows. Section II presents the HT method: first, the Hunt filter is briefly recalled; then, we derive the approximate pdf of the Hunt filter output from which we determine the threshold that minimizes the probability of errors. In Section III, based on an interpretation of the HT method as a maximum *a posteriori* (MAP) estimator, we propose to estimate the hyperparameters using a joint MAP (JMAP) approach [7]. In Section IV, the performances of the HT method are evaluated and compared to those of the SMLR of [1]. Finally, Section V gives the conclusions and perspectives of this letter.

II. HT METHOD

A. Hunt Filter

The Hunt filter [6] results from the Phillips and Twomey criterion minimization, which corresponds to the discrete time formulation of the Tikhonov criterion

$$\mathcal{J}_{\text{PT}} = (\mathbf{y} - \mathbf{H}\mathbf{x})^T (\mathbf{y} - \mathbf{H}\mathbf{x}) + \alpha \mathbf{x}^T \mathbf{D}^T \mathbf{D} \mathbf{x}$$

where $\mathbf{D} = \mathbf{I}_N$ so as to favor the restoration of zero-value signal.

Minimizing \mathcal{J}_{PT} yields a first estimation $\tilde{\mathbf{x}}$ of \mathbf{x}

$$\tilde{\mathbf{x}} = \mathbf{G}\mathbf{y} \quad \text{where} \quad \mathbf{G} = (\mathbf{H}^T \mathbf{H} + \alpha \mathbf{I}_N)^{-1} \mathbf{H}^T. \quad (1)$$

The matrix \mathbf{H} being supposed circulant, its (discrete) Fourier transform is a diagonal matrix. Thus, the algorithm may be efficiently implemented using the fast Fourier transform, $\tilde{\mathbf{x}}$ being then obtained by an inverse (discrete) Fourier transform. This results in the so-called Hunt filter.

A Bayesian interpretation [5] shows that the Hunt filter is equivalent to a MAP approach with a Gaussian prior: it yields $\alpha = r_n/r_x$, where r_x is the variance of the Gaussian signal to restore. In the sequel, we will also take $\alpha = r_n/r_x$ with r_x

the variance of the pulses of the BG sequence to favor the pulse amplitude estimation.

B. Approximate PDF of $\tilde{\mathbf{x}}$

The pdf $p(\tilde{\mathbf{x}})$ of $\tilde{\mathbf{x}} = \mathbf{G}\mathbf{y} = \mathbf{G}\mathbf{H}\mathbf{x} + \mathbf{G}\mathbf{n}$ is needed to calculate the threshold t (Section II-C).

As $(\mathbf{G}\mathbf{n})_k = \sum_{i=1}^N \mathbf{G}_{k,i} \mathbf{n}_i$, we have

$$p((\mathbf{G}\mathbf{n})_k) = p\left(\sum_{i=1}^N \mathbf{G}_{k,i} \mathbf{n}_i\right) = \prod_{i=\{1,\dots,N\}}^* p(\mathbf{G}_{k,i} \mathbf{n}_i).$$

\prod^* is introduced to improve the readability of the letter

$$\prod_{i=\{1,\dots,N\}}^* f_i(x) = [f_1 * \dots * f_N](x).$$

Since $p(\mathbf{n}_k) = \mathcal{N}(0, r_n)$ and $\mathcal{N}(0, \sigma_1^2) * \mathcal{N}(0, \sigma_2^2) = \mathcal{N}(0, \sigma_1^2 + \sigma_2^2)$, we get

$$p((\mathbf{G}\mathbf{n})_k) = \mathcal{N}\left(0, r_n \sum_{i=1}^N \mathbf{G}_{k,i}^2\right). \quad (2)$$

Similarly, since $p(\mathbf{x}_k) = \lambda \mathcal{N}(0, r_x) + (1 - \lambda) \delta(\mathbf{x}_k)$ and $(\mathbf{G}\mathbf{H}\mathbf{x})_k = \sum_{i=1}^N (\mathbf{G}\mathbf{H})_{k,i} \mathbf{x}_i$, we have

$$\begin{aligned} p((\mathbf{G}\mathbf{H}\mathbf{x})_k) &= \prod_{i=\{1,\dots,N\}}^* p((\mathbf{G}\mathbf{H})_{k,i} \mathbf{x}_i) \\ &= \prod_{i=\{1,\dots,N\}}^* [\lambda \mathcal{N}(0, r_x (\mathbf{G}\mathbf{H})_{k,i}^2) + (1 - \lambda) \delta(\mathbf{x}_i)]. \end{aligned}$$

To simplify the problem, we now need two assumptions.

- 1) We suppose that the off-diagonal terms of $\mathbf{G}\mathbf{H}$ are negligible as compared to the diagonal terms. This is nothing but considering the values $(\mathbf{G}\mathbf{H}\mathbf{x})_k$ decorrelated. We note that the higher the SNR is (i.e., α decreases), the more valid this approximation is. Indeed, when α decreases, $\mathbf{G}\mathbf{H}$ tends to \mathbf{I}_N . With this assumption, we have $r_x (\mathbf{G}\mathbf{H})_{k,i} = 0 \forall k \neq i$. The validity of this approximation depends both on the SNR and the impulse response. Typically, for the impulse response used in this letter (see Section IV), the approximation is very realistic for a SNR greater than 15 dB. But if the frequency contents of the impulse response becomes more concentrated in low frequencies, the SNR should be greater;
- 2) We suppose that the diagonal terms of $\mathbf{G}\mathbf{G}^T$ and $\mathbf{G}\mathbf{H}\mathbf{H}^T \mathbf{G}^T$ are constant, i.e., we neglect the boundary effects.

Thus, considering that a Gaussian with zero variance is a Dirac mass, we get

$$p((\mathbf{G}\mathbf{H}\mathbf{x})_k) = \lambda \mathcal{N}(0, r_x (\mathbf{G}\mathbf{H})_{k,k}^2) + (1 - \lambda) \delta(\mathbf{x}_k). \quad (3)$$

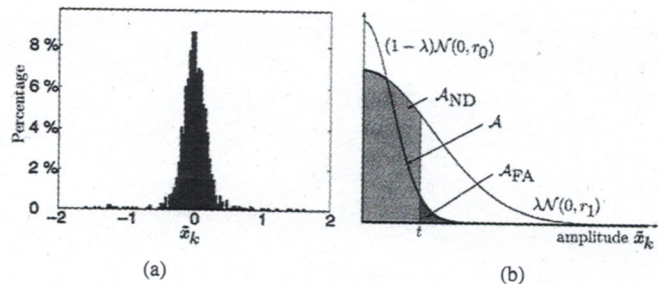


Fig. 1. PDF of $\tilde{\mathbf{x}}$. (a) Histogram of $\tilde{\mathbf{x}}$ for 1024 samples. (b) Areas of false alarms and nondetections.

From (2) and (3), we obtain

$$\begin{aligned} p(\tilde{\mathbf{x}}_k) &= p((\mathbf{G}\mathbf{H}\mathbf{x})_k) * p((\mathbf{G}\mathbf{n})_k) \\ &= \lambda \mathcal{N}\left(0, r_x (\mathbf{G}\mathbf{H})_{k,k}^2 + r_n \sum_{i=1}^N \mathbf{G}_{k,i}^2\right) \\ &\quad + (1 - \lambda) \mathcal{N}\left(0, r_n \sum_{i=1}^N \mathbf{G}_{k,i}^2\right). \end{aligned}$$

As $r_n \sum_{i=1}^N \mathbf{G}_{k,i}^2 = r_n (\mathbf{G}\mathbf{G}^T)_{k,k}$ and $r_x (\mathbf{G}\mathbf{H})_{k,k}^2 = r_x \sum_{i=1}^N (\mathbf{G}\mathbf{H})_{k,i}^2 = r_x (\mathbf{G}\mathbf{H}\mathbf{H}^T \mathbf{G}^T)_{k,k}$, it follows that the pdf is a Gaussian mixture centered on zero

$$p(\tilde{\mathbf{x}}_k) = \lambda \mathcal{N}(0, r_1) + (1 - \lambda) \mathcal{N}(0, r_0) \quad (4)$$

with

$$r_1 = r_x (\mathbf{G}\mathbf{H}\mathbf{H}^T \mathbf{G}^T)_{k,k} + r_n (\mathbf{G}\mathbf{G}^T)_{k,k} \quad (5)$$

$$r_0 = r_n (\mathbf{G}\mathbf{G}^T)_{k,k}, \quad r_0 < r_1. \quad (6)$$

The histogram of $\tilde{\mathbf{x}}$ obtained on a 1024-sample signal [Fig. 1(a)] clearly shows that the Gaussian mixture is a good model for $p(\tilde{\mathbf{x}}_k)$.

C. Thresholding

The thresholding separates pulses [term $\lambda \mathcal{N}(0, r_1)$] from noise [term $(1 - \lambda) \mathcal{N}(0, r_0)$]. The estimated signal may be expressed as $\hat{\mathbf{x}}_k = \tilde{\mathbf{x}}_k$ if $|\tilde{\mathbf{x}}_k| > t$, $\hat{\mathbf{x}}_k = 0$ otherwise. Performing such a (hard) thresholding, one can make two kinds of errors. *False alarms* correspond to wrongly detected pulses, while *nondetections* correspond to missed pulses. The threshold t is chosen to minimize the probability of error $\{p_{FA} + p_{ND}\}$ where p_{FA} and p_{ND} are the probabilities of false alarm and nondetection. Fig. 1(b) shows the positive part of the two Gaussians. The area \mathcal{A}_{FA} represents half the probability p_{FA} and \mathcal{A}_{ND} half the probability p_{ND} . It appears that $\mathcal{A}_{ND} + \mathcal{A}_{FA}$ is minimal and equal to \mathcal{A} when t corresponds to the intersection point of the two Gaussians, which gives

$$t = \sqrt{\ln\left(\frac{\lambda}{1 - \lambda} \sqrt{\frac{r_0}{r_1}}\right) \frac{2r_1 r_0}{r_0 - r_1}}. \quad (7)$$

III. HYPERPARAMETER ESTIMATION

To address the hyperparameter estimation problem, we first need to interpret the HT method in a Bayesian framework. In [8] and [9], it is shown that the MAP estimator remains unchanged

for data varying in a neighborhood if and only if the prior pdf is nonsmooth (i.e., its derivative is discontinuous) in this neighborhood. This results in a thresholding effect of the MAP estimator under nonsmooth priors. Our prior being BG, the MAP estimation implies a thresholding, that allows us to use the HT method as a MAP estimator. The BG estimation $\hat{\mathbf{x}}$ may be interpreted as the pointwise multiplication: $\forall k, \hat{\mathbf{x}}_k = \tilde{\mathbf{x}}_k \hat{q}_k$, $\tilde{\mathbf{x}}$ corresponding to the Hunt filter output and \hat{q} being an estimate of the i.i.d. Bernoulli sequence with parameter λ , controlling the occurrence of pulses in \mathbf{x} . So, we consider the MAP estimation of (\mathbf{x}, \mathbf{q}) , which maximizes the following joint posterior pdf:

$$p(\mathbf{x}, \mathbf{q} | \mathbf{y}, \theta) \propto p(\mathbf{y} | \mathbf{x}, \mathbf{q}, \theta) p(\mathbf{x} | \mathbf{q}, \theta) p(\mathbf{q} | \theta) \quad (8)$$

where

- $p(\mathbf{y} | \mathbf{x}, \mathbf{q}, \theta) = p(\mathbf{y} | \mathbf{x}, \theta) = \mathcal{N}(\mathbf{H}\mathbf{x}, r_n \mathbf{I}_N)$ because $p(\mathbf{n} | \theta) = \mathcal{N}(\mathbf{0}, r_n \mathbf{I}_N)$ and $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n}$;
- $p(\mathbf{x} | \mathbf{q}, \theta) = \mathcal{N}(\mathbf{0}, r_x \mathbf{Q})$ where $\mathbf{Q} = \text{diag}\{\mathbf{q}\}$;
- $p(\mathbf{q} | \theta) = \lambda^{N_q} (1 - \lambda)^{N - N_q}$, N_q being the number of nonzero samples of \mathbf{q} .

To estimate the three hyperparameters gathered into $\theta = \{\lambda, r_x, r_n\}$, we consider the following JMAP [7]:

$$\begin{aligned} (\hat{\mathbf{x}}, \hat{\mathbf{q}}, \hat{\theta}) &= \arg \max_{(\mathbf{x}, \mathbf{q}, \theta)} p(\mathbf{x}, \mathbf{q}, \theta | \mathbf{y}) \\ &= \arg \max_{(\mathbf{x}, \mathbf{q}, \theta)} p(\mathbf{y} | \mathbf{x}, \mathbf{q}, \theta) p(\mathbf{x} | \mathbf{q}, \theta) p(\mathbf{q} | \theta) \end{aligned}$$

where $p(\theta)$ does not appear because it is considered as uniform (no *a priori* on the hyperparameters: it is in fact an ML estimator). This optimization problem is solved using an iterative procedure

$$\begin{cases} (\hat{\mathbf{x}}^{(i)}, \hat{\mathbf{q}}^{(i)}) = \arg \max_{(\mathbf{x}, \mathbf{q})} p(\mathbf{x}, \mathbf{q}, \lambda^{(i-1)}, r_x^{(i-1)}, r_n^{(i-1)} | \mathbf{y}) \\ (\hat{\lambda}^{(i)}, \hat{r}_x^{(i)}, \hat{r}_n^{(i)}) = \arg \max_{(\lambda, r_x, r_n)} p(\mathbf{x}^{(i)}, \mathbf{q}^{(i)}, \lambda, r_x, r_n | \mathbf{y}). \end{cases}$$

The first optimization problem is approximatively solved using the HT method, and the hyperparameters are estimated using (8), with the assumption that \mathbf{x} is a BG signal

$$\begin{aligned} \hat{\lambda} &= \arg \max_{\lambda} p(\mathbf{q} | \lambda) = \frac{N_q}{N} \\ \hat{r}_x &= \arg \max_{r_x} p(\mathbf{x} | \mathbf{q}, r_x) = \frac{1}{N_q} \|\hat{\mathbf{x}}\|^2 \\ \hat{r}_n &= \arg \max_{r_n} p(\mathbf{y} | \mathbf{x}, \mathbf{q}, r_n) = \frac{1}{N} \|\mathbf{y} - \mathbf{H}\hat{\mathbf{x}}\|^2. \end{aligned}$$

Rather than using (5) and (6) to estimate r_0 and r_1 , we propose to estimate them directly from $\hat{\mathbf{x}}$ as

$$(\hat{r}_0, \hat{r}_1) = \arg \max_{(r_0, r_1)} p(\hat{\mathbf{x}} | \mathbf{q}, r_0, r_1).$$

From (4), we get for all k

$$\begin{aligned} p(\hat{\mathbf{x}}_k | \mathbf{q}, r_0, r_1) &= \mathbf{q}_k \mathcal{N}(0, r_1) + (1 - \mathbf{q}_k) \mathcal{N}(0, r_0) \\ &= \mathcal{N}(0, r_1 \mathbf{q}_k + r_0(1 - \mathbf{q}_k)). \end{aligned}$$

Then

$$p(\hat{\mathbf{x}} | \mathbf{q}, r_0, r_1) = \mathcal{N}(\mathbf{0}, r_1 \mathbf{Q} + r_0(\mathbf{I}_N - \mathbf{Q})).$$

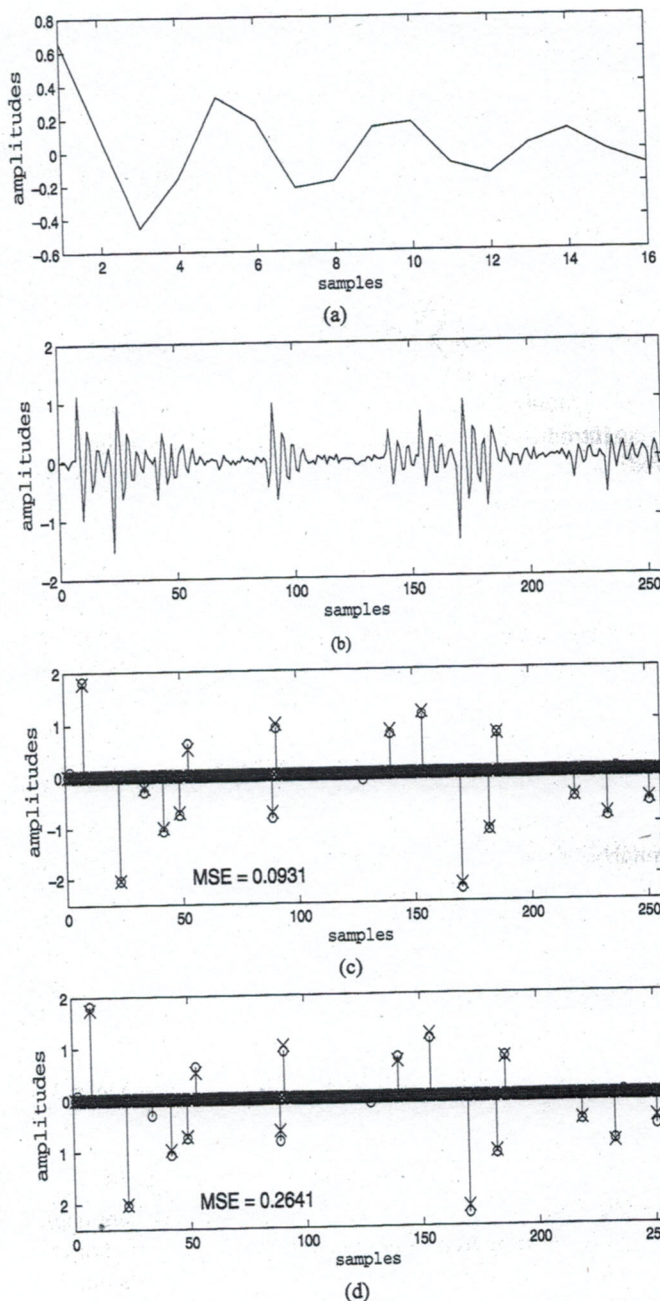


Fig. 2. SMLR and HT estimates. \circ represents real pulses, and \times estimated ones. MSE: mean square error on the detected pulses. (a) Impulse response (h). (b) Data (y). (c) SMLR estimate. (d) HT estimate.

Straightforward calculations yield the following explicit expressions:

$$\hat{r}_1 = \frac{1}{N_q} \sum_{k=1}^N \hat{\mathbf{x}}_k^2 \hat{q}_k, \quad \hat{r}_0 = \frac{1}{N - N_q} \sum_{k=1}^N \hat{\mathbf{x}}_k^2 (1 - \hat{q}_k).$$

This approach has been chosen because it leads to a faster algorithm and gives better results than those obtained by using (5) and (6).

λ , r_x , and r_n are initialized to arbitrary values, while r_0 and r_1 are initialized using (5) and (6). The convergence test is made by comparing the signal and hyperparameter values from one iteration to the next. The procedure stops when they differ no

TABLE I
COMPUTATION TIME, PERCENTAGE OF DETECTED PULSES (DP), AND
FALSE ALARMS (FA) FOR SEVERAL SNR

SNR =	SMLR	23.39 s	92.0 % DP	2.0 % FA
20 dB	HT	0.16 s	87.0 % DP	5.0 % FA
SNR =	SMLR	26.94 s	88.17 % DP	2.15 % FA
15 dB	HT	0.19 s	81.72 % DP	4.3 % FA
SNR =	SMLR	37.59 s	69.07 % DP	3.09 % FA
10 dB	HT	0.23 s	61.85 % DP	8.24 % FA
SNR =	SMLR	105.64 s	42.9 % DP	3.7 % FA
5 dB	HT	0.31 s	40.2 % DP	1.9 % FA

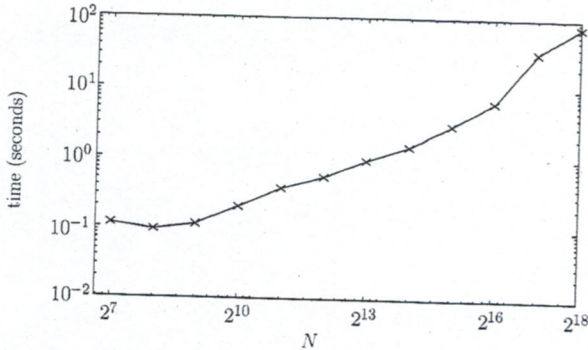


Fig. 3. Mean computation time.

more than 1%. Extensive simulations have shown the very satisfactory behavior of the procedure.

IV. SIMULATIONS RESULTS

Some simulations were carried out to assess the performances of the HT method and to compare them to those of the SMLR of [1]. The two methods were implemented in Matlab 6.5 on a Pentium III at 800 MHz. Fig. 2 shows the results achieved by the two methods on a 256-sample signal with a 16-sample impulse response and the following real hyperparameters: $\lambda = 0.08$, $\alpha = 1.5$, SNR = 15 dB. Comparing the results, it appears that the SMLR performs slightly better than the HT method. Table I gives the statistics obtained on ten simulations. They confirm the superiority of the SMLR at the price of a very high computation time as compared to the HT method. We also note that

the performances of the two methods become comparable as the SNR increases. Fig. 3 shows the mean computation time evolution of the HT method as a function of N , confirming that the computational burden remains reasonable, even for very large signals.

V. CONCLUSION

The HT method is a sparse spike train deconvolution consisting in coupling the Hunt filter with a thresholding. We show that when the signal x is BG, a Gaussian mixture is a good model for the pdf of the Hunt filter output. From this result, we derive the threshold that minimizes the probability of errors. We then propose to use the HT method as a MAP estimator and to jointly estimate the hyperparameters of the problem. Simulations show that the method performs well at a very low computation time, making this approach very well-suited to process very large signals and to be implemented on embedded systems. Future work will be directed to find which global criterion is minimized.

REFERENCES

- [1] F. Champagnat, Y. Goussard, and J. Idier, "Unsupervised deconvolution of sparse spike trains using stochastic approximation," *IEEE Trans. Signal Processing*, vol. 44, pp. 2988–2998, Dec. 1996.
- [2] J. J. Kormylo and J. M. Mendel, "Maximum likelihood detection and estimation of Bernoulli–Gaussian processes," *IEEE Trans. Inform. Theory*, vol. IT-28, pp. 482–488, May 1982.
- [3] V. Mazet, D. Brie, and C. Caironi, "Déconvolution impulsionnelle par filtre de Hunt et seuillage," in *Proc. GRETSI*, Sept. 2003.
- [4] K. F. Kaarensen, "Deconvolution of sparse spike trains by iterated window maximization," *IEEE Trans. Signal Processing*, vol. 45, pp. 1173–1183, May 1997.
- [5] J. Idier, *Approche bayésienne pour les problèmes inverses*. Paris, France: Hermès Science, 2001.
- [6] B. R. Hunt, "The inverse problem of radiography," *Math. Biosci.*, vol. 8, pp. 161–179, 1970.
- [7] A. Mohammad-Djafari, "Joint estimation of parameters and hyperparameters in a Bayesian approach of solving inverse problems," in *Proc. ICASSP*, Munich, Germany, Apr. 1997, pp. 2837–2840.
- [8] P. Moulin and J. Liu, "Analysis of multiresolution image denoising schemes using generalized Gaussian and complexity priors," *IEEE Trans. Inform. Theory*, vol. 45, pp. 909–919, Apr. 1999.
- [9] M. Nikolova, "Local strong homogeneity of a regularized estimator," *SIAM J. Appl. Math.*, vol. 61, no. 2, pp. 633–658, 2000.