

# Kullback Proximal Algorithms for Maximum Likelihood Estimation

Stéphane Chrétien, Alfred O. Hero

► **To cite this version:**

Stéphane Chrétien, Alfred O. Hero. Kullback Proximal Algorithms for Maximum Likelihood Estimation. [Research Report] RR-3756, INRIA. 1999. <inria-00072906>

**HAL Id: inria-00072906**

**<https://hal.inria.fr/inria-00072906>**

Submitted on 24 May 2006

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

***Kullback Proximal Algorithms for Maximum  
Likelihood Estimation***

Stéphane Chrétien, Alfred O. Hero

**N° 3756**

Août 1999

———— THÈME 4 ————



*Rapport  
de recherche*



# Kullback Proximal Algorithms for Maximum Likelihood Estimation

Stéphane Chrétien, Alfred O. Hero\*

Thème 4 — Simulation et optimisation  
de systèmes complexes  
Projet is2

Rapport de recherche n° 3756 — Août 1999 — 15 pages

**Abstract:** In this paper, we study the convergence of a new class of fast and stable sequential optimization methods for computing maximum likelihood estimates. These methods are based on a *proximal point algorithm* implemented with a Kullback-type proximal function. When the proximal regularization parameter is set to unity one obtains the classical expectation maximization (EM) algorithm. For other values of the regularization parameter, relaxed versions of EM are obtained which can have much faster convergence. In particular, if the regularization parameter vanishes at infinity, a superlinearly converging algorithm is obtained. We present an implementation of the algorithm using the *trust region* update strategy. For illustration the method is applied to a non-quadratic inverse problem with Poisson distributed data.

**Key-words:** accelerated EM algorithm, Kullback regularization, proximal point iterations, superlinear convergence, trust region methods.

(Résumé : *tsvp*)

\* Dept. of Electrical Engineering and Computer Science, The University of Michigan, 1301 Beal Avenue, Ann Arbor, Michigan 48109-2122, USA, Email hero@eecs.umich.edu

# Méthodes Kullback-proximales pour l'estimation au sens du maximum de vraisemblance

**Résumé :** Dans cet article, nous étudions la convergence d'une nouvelle classe de méthodes rapides et stables d'optimisation servant à obtenir des estimateurs au sens du maximum de vraisemblance. Ces méthodes sont fondées sur l'algorithme du point proximal utilisant une régularisation du type divergence de Kullback. Quand le paramètre de relaxation est l'unité, on obtient l'algorithme EM classique comme cas particulier. Pour d'autres valeurs du paramètre, des versions relaxées de EM sont obtenues, pouvant jouir d'une convergence beaucoup plus rapide. En particulier, si le paramètre de relaxation converge vers zéro, on obtient un algorithme à convergence superlinéaire. Nous présentons aussi une implémentation par régions de confiance. Le comportement de ces méthodes est illustré sur un exemple, un problème inverse avec données poissonniennes.

**Mots-clé :** accélération de l'algorithme EM, régularisation, divergence de Kullback, algorithme du point proximal, convergence superlinéaire, régions de confiance.

## 1 Introduction

Iterative solutions to the maximum likelihood (ML) estimation problem are of interest when direct closed form solution is infeasible. Among the most stable iterative strategies to ML is the popular expectation maximization (EM) algorithm [1]. This algorithm has the attractive property of monotonicity which guarantees that the likelihood function increases with each iteration. The convergence properties of the EM algorithm and its variants have been extensively studied in the literature; see [2] and [3] for instance. It is well known that under strong concavity assumptions the EM algorithm converges linearly towards the ML estimator  $\theta_{ML}$ . However, in practice the EM algorithm suffers from slow convergence in late iterations. Efforts to improve the asymptotic convergence of the EM algorithm include: Aitken's acceleration [4], overrelaxation [5], conjugate gradient [6] [7], Newton methods [8] [9], quasi-Newton methods [10], ordered subsets EM [11] and stochastic EM [12]. Unfortunately, these methods do not automatically guarantee the monotone increasing likelihood property of standard EM, which frequently requires additional monitoring for instability [13].

The main goal of this paper is to recast the EM algorithm into a more general framework of monotone algorithms having the potential for accelerated convergence. For this purpose, the EM algorithm is identified as a particular instance of a proximal point algorithm using Kullback regularization. The proximal point algorithm, first introduced by Rockafellar [14] and Martinet [15], is a state of the art procedure in optimization. In particular, proximal approaches have led to many high performance numerical algorithms; e.g. bundle methods for nonsmooth problems [16] and multiplier methods for constrained optimization [17]. A key motivation for the proximal point algorithm is that an iteration-dependent penalty can be introduced to obtain superlinear convergence rates in the case of quadratic regularization [14]. In this paper, this idea is used to obtain relaxed versions of EM algorithm with superlinear asymptotic convergence rates.

The outline of the paper is the following. In Section 2 we provide a brief review of key elements of the classical EM algorithm. In Section 3, we establish a relationship between the EM algorithm and the proximal point algorithm. In section 4, we present the general Kullback proximal point algorithm and we establish global and superlinear convergence to the maximum likelihood estimator. In section 5, we study second order approximations of the Kullback proximal point iteration using trust region updating. The trust region strategy is introduced in order to obtain global convergence and monotonic behavior of the approximate scheme. Finally, in Section 6 we present numerical comparisons for a Poisson inverse problem.

Our notations are standard for the most part.  $\nabla_{10}I(\theta, \bar{\theta})$  (resp.  $\nabla_{10}^2I(\theta, \bar{\theta})$ ) denotes the gradient (resp. the hessian matrix) of  $I(\theta, \bar{\theta})$  in the first variable.

## 2 Background

The problem of maximum likelihood (ML) estimation consists of finding a solution to

$$\theta_{ML} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} l_y(\theta), \quad (1)$$

where  $y$  is an observed sample of a random variable  $Y$  defined on a sample space  $\mathcal{Y}$  and  $l_y(\theta)$  is the log-likelihood function defined by

$$l_y(\theta) = \log g(y; \theta), \quad (2)$$

where  $g(y; \theta)$  denotes the density of  $Y$  at  $y$  parametrized by a deterministic vector  $\theta$  in  $\mathbb{R}^p$ . One of the most popular methods for solving ML estimation problems is the Expectation Maximization (EM) algorithm described in Dempster, Laird, and Rubin [1] which we describe as follows.

A more informative data space  $\mathcal{X}$  is introduced. A random variable  $X$  is defined on  $\mathcal{X}$  that with density  $f(x; \theta)$  parametrized by  $\theta$ . The data  $X$  is more informative than the actual data  $Y$  in the sense that  $Y$  is a compression of  $X$ , i.e. there exists a non-invertible transformation  $h$  such that  $Y = h(X)$ . It would therefore be advantageous to replace the ML estimation problem (1) by

$$\theta_{ML} = \max_{\theta \in \mathbb{R}^p} l_x(\theta), \quad (3)$$

with  $l_x(\theta) = \log f(x; \theta)$ . Since  $y = h(x)$  the density  $g$  of  $Y$  is related to the density  $f$  of  $X$  through

$$g(y; \theta) = \int_{h^{-1}(\{y\})} f(x; \theta) d\mu(x) \quad (4)$$

for an appropriate measure  $\mu$  on  $\mathcal{X}$ . Under condition (4) and for a given observed sample  $y$ , for any  $x$  in  $h^{-1}(\{y\})$  the solutions of (3) are solutions of the original ML estimation problem (1). In this setting, the data  $y$  are called *incomplete data* whereas the data  $x$  are called *complete data*.

It remains to deal with the fact that the complete data  $x$  corresponding to a given observed sample  $y$  are unknown. Therefore, the complete data likelihood function  $l_x(\theta)$  can only be estimated. Given the observed data  $y$  and a previous estimate of  $\theta$  denoted  $\bar{\theta}$ , the following minimum mean square error estimator (MMSE) of the quantity  $l_x(\theta)$  is natural

$$Q(\theta, \bar{\theta}) = \mathbf{E}[\log f(x; \theta) | y; \bar{\theta}],$$

where, for any integrable function  $F(x)$  on  $\mathcal{X}$ , we have defined the conditional expectation

$$\mathbf{E}[F(x) | y; \bar{\theta}] = \int_{h^{-1}(\{y\})} F(x) k(x | y; \bar{\theta}) d\mu(x)$$

and  $k(x | y; \bar{\theta})$  is the conditional density function given  $y$

$$k(x | y; \bar{\theta}) = \frac{f(x; \bar{\theta})}{g(y; \bar{\theta})}. \quad (5)$$

The EM algorithm generates a sequence of approximations to the solution (3) starting from an initial guess  $\theta^0$  of  $\theta_{ML}$  and defined by

$$\mathbf{Compute} \quad Q(\theta, \theta^k) = \mathbf{E}[\log f(x; \theta) | y; \theta^k] \quad \mathbf{E \ Step}$$

$$\theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} Q(\theta, \theta^k) \quad \mathbf{M \ Step}$$

The key to understanding the convergence of the EM algorithm is the following decomposition of the log-likelihood function

$$l_y(\theta) = Q(\theta, \bar{\theta}) + H(\theta, \bar{\theta}) \quad (6)$$

where

$$H(\theta, \bar{\theta}) = -\mathbf{E}[k(x | y; \theta) | y; \bar{\theta}].$$

It follows from elementary application of Jensen's inequality to the log function that

$$H(\theta, \bar{\theta}) \geq H(\theta, \theta) \geq 0, \quad \forall \theta, \bar{\theta} \in \mathbb{R}^p. \quad (7)$$

Observe from (6) and (7) that the  $\theta$  function  $\{Q(\theta, \theta^k)\}$  is a lower bound on the log likelihood function  $l_y(\theta)$ . This fact is necessary to ensure monotonicity of the algorithm. Specifically, using the M-step defining relation

$$Q(\theta^{k+1}, \theta^k) \geq Q(\theta^k, \theta^k), \quad (8)$$

one obtains

$$l_y(\theta^{k+1}) - l_y(\theta^k) \geq Q(\theta^{k+1}, \theta^k) - Q(\theta^k, \theta^k) \quad (9)$$

$$+ H(\theta^{k+1}, \theta^k) - H(\theta^k, \theta^k). \quad (10)$$

Hence, using (8) and (7)

$$l_y(\theta^{k+1}) \geq l_y(\theta^k).$$

This is the well known monotonicity property of the EM algorithm.

Note that if the function  $H(\theta, \bar{\theta})$  in (6) were scaled by an arbitrary positive factor  $\beta$  the function  $Q(\theta, \bar{\theta})$  would remain a lower bound on  $l_y(\theta)$  and monotonicity of the algorithm would be preserved. As will be shown below, if  $\beta$  is allowed to vary with iteration in a suitable manner one obtains a monotone, superlinearly convergent generalization of the EM algorithm.

### 3 Proximal methods and the EM algorithm

In this section, we present the proximal point algorithm of Rockafellar and Martinet. We then demonstrate that EM is a particular instance of a proximal point method implemented with a Kullback-type penalty.

#### 3.1 The proximal point algorithm

Consider the general problem of maximizing a concave function  $\Phi(\theta)$ . The proximal point algorithm is an iterative procedure which can be written

$$\theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \left\{ \Phi(\theta) - \frac{\beta_k}{2} \|\theta - \theta^k\|^2 \right\}. \quad (11)$$

The quadratic penalty  $\|\theta - \theta^k\|^2$  is relaxed using a sequence of positive parameters  $\{\beta_k\}$ . In [14], Rockafellar showed that superlinear convergence of this method is obtained when the sequence  $\{\beta_k\}$  converges towards zero. In numerical implementations of proximal point the function  $\Phi(\theta)$  is generally replaced by a piecewise linear model [16].

#### 3.2 Proximal interpretation of the EM algorithm

We now turn to the relation between the EM algorithm and the proximal point framework. For this purpose, we will need to consider a particular Kullback information measure. Assume that the family  $\{k(x|y; \theta)\}_{\theta \in \mathbb{R}^p}$  is regular in the sense of Ibragimov and Khasminskii [18], in particular  $k(x|y; \theta)\mu(x)$  and  $k(x|y; \bar{\theta})\mu(x)$  are mutually absolutely continuous for any  $\theta$  and  $\bar{\theta}$  in  $\mathbb{R}^p$ . Then the Radon-Nikodym derivative  $\frac{k(x|y; \bar{\theta})}{k(x|y; \theta)}$  exists for all  $\theta, \bar{\theta}$  and we can define the following Kullback measure.

$$I(\theta, \bar{\theta}|y) = \mathbb{E} \left[ \log \frac{k(x|y, \bar{\theta})}{k(x|y; \theta)} \middle| y; \bar{\theta} \right]. \quad (12)$$

**Proposition 1** *The EM algorithm is a proximal point algorithm with Kullback-type penalty (12) of the form*

$$\theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \{ l_y(\theta) - I(\theta, \bar{\theta}|y) \} \quad (13)$$

**Proof 1** *The key to making the connection with the proximal point algorithm is the following representation of the M step:*

$$\theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \left\{ \log g(y; \theta) + \mathbb{E} \left[ \log \frac{f(x; \theta)}{g(y; \theta)} \middle| y; \theta^k \right] \right\}.$$

*This equation is equivalent to*

$$\begin{aligned} \theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \left\{ \log g(y; \theta) + \mathbb{E} \left[ \log \frac{f(x; \theta)}{g(y; \theta)} \middle| y; \theta^k \right] \right. \\ \left. - \mathbb{E} \left[ \log \frac{f(x; \theta^k)}{g(y; \theta^k)} \middle| y; \theta^k \right] \right\} \end{aligned}$$

*since the additional term does not modify the maximization problem. Recalling that  $k(x|y; \theta) = \frac{f(x; \theta)}{g(y; \theta)}$ ,*

$$\begin{aligned} \theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \left\{ \log g(y; \theta) + \mathbb{E} [k(x|y; \theta) | y; \theta^k] \right. \\ \left. - \mathbb{E} [k(x|y; \theta^k) | y; \theta^k] \right\}. \end{aligned}$$

*We finally obtain*

$$\theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \left\{ \log g(y; \theta) + \mathbb{E} \left[ \frac{k(x|y; \theta)}{k(x|y; \theta^k)} \middle| y; \theta^k \right] \right\}$$

*which concludes the proof.*



## 4 Kullback Generalization of Proximal Point Algorithm

The proximal point interpretation of EM in the previous section suggests that EM can be generalized by replacing the quadratic penalty in the standard proximal point recursion (11) with the Kullback penalty (12). Other non-quadratic generalizations of proximal point have been proposed by Teboulle and others (see [19] and references therein). Moreover, the parallel with the proximal point method also suggests that our approach applies to a large class of (twice differentiable) objective functions, including penalized Likelihood functions. In this section we introduce the Kullback generalization of proximal point, called the Kullback proximal point (KPP) algorithm, and establish its convergence properties.

### 4.1 Algorithm Definition

Let  $\{k(x; \theta)\}_{\theta \in \mathbb{R}^p}$  be any family of densities supported on a set  $S \subset \mathcal{X}$ , parametrized by  $\theta$  and regular in the sense of Ibragimov and Khasminskii [18]. Define the Kullback information measure

$$I(\theta, \bar{\theta}) = \int_S \log \frac{k(x; \theta)}{k(x; \bar{\theta})} k(x; \bar{\theta}) d\mu(x). \quad (14)$$

Notice that for the choice  $k(x; \theta) = k(x|y; \theta)$  and for  $S = h^{-1}(\{y\})$ , we obtain  $I(\theta, \bar{\theta}) = I(\theta, \bar{\theta}|y)$  as defined by (12) in the context of the EM algorithm. We make the following assumptions, where  $\Lambda_M$  denotes the greatest eigenvalue of the matrix  $M$  and  $\lambda_M$  denotes the smallest.

**Assumptions 1** *We assume the following:*

- (i)  $l_y(\theta)$  is twice continuously differentiable on  $\mathbb{R}^p$  and  $I(\theta, \bar{\theta})$  is twice continuously differentiable in  $(\theta, \bar{\theta})$  in  $\mathbb{R}^p \times \mathbb{R}^p$ .
- (ii)  $\lim_{\|\theta\| \rightarrow \infty} l_y(\theta) = -\infty$ .
- (iii)  $l_y(\theta) < \infty$  and  $\lambda_{\nabla^2 l_y(\theta)} \leq \Lambda_{\nabla^2 l_y(\theta)} < 0$  on every bounded  $\theta$ -set
- (iv) for any  $\bar{\theta}$  in  $\mathbb{R}^p$ ,  $I(\theta, \bar{\theta}) < \infty$  and  $0 < \lambda_{\nabla_{\theta}^2 I(\theta, \bar{\theta})} \leq \Lambda_{\nabla_{\theta}^2 I(\theta, \bar{\theta})}$  on every bounded  $\theta$ -set.

These assumptions ensure smoothness of  $l_y(\theta)$  and  $I(\theta, \bar{\theta})$  and their first two derivatives in  $\theta$ . Assumption 1.iii also implies strong concavity of  $l_y(\theta)$ . Assumption 1.iv implies that  $I(\theta, \bar{\theta})$  is strictly convex and that the parameter  $\theta$  is strongly identifiable in the family of densities  $k(x; \theta)$  (see proof of Lemma 1 below). Note that the above assumptions are not the minimum set possible, e.g. that  $l_y(\theta)$  and  $I(\theta, \bar{\theta})$  are upper bounded follows from continuity, Assumption 1.ii and the property  $I(\theta, \bar{\theta}) \geq I(\bar{\theta}, \bar{\theta}) = 0$ , respectively.

With the above assumptions we are now prepared to introduce the general algorithm.

**Definition 1** *Let  $g(y; \theta)$  and  $k(x; \theta)$  be such that  $l_y(\theta)$  and  $I(\theta, \bar{\theta})$  satisfy Assumptions 1. Let  $\{\beta_k\}$  be a sequence of positive relaxation parameters. Then, the following recurrence will be called the Kullback proximal point (KPP) algorithm*

$$\theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \{l_y(\theta) - \beta_k I(\theta, \theta^k)\}. \quad (15)$$

The KPP algorithm is well defined since the maximum in (15) is always achieved in a bounded set. Monotonicity is guaranteed by this procedure as proved in the following proposition.

**Proposition 2** *The log-likelihood sequence  $\{l_y(\theta^k)\}$  is monotone non-decreasing and satisfies*

$$l_y(\theta^{k+1}) - l_y(\theta^k) \geq \beta_k I(\theta^{k+1}, \theta^k), \quad (16)$$

**Proof 2** *From iteration (15), we have*

$$l_y(\theta^{k+1}) - l_y(\theta^k) \geq \beta_k I(\theta^{k+1}, \theta^k) - \beta_k I(\theta^k, \theta^k).$$

*Since  $I(\theta^k, \theta^k) = 0$  and  $I(\theta^{k+1}, \theta^k) \geq 0$ , we deduce (16) and that  $\{l_y(\theta^k)\}$  is non-decreasing.*

We next turn to convergence of the KPP iterates  $\{\theta^k\}$ .

## 4.2 Convergence

First we characterize the fixed points of the KPP algorithm.

**Proposition 3** *Given  $\beta_k = \beta > 0$ ,  $k = 1, 2, \dots$ , the fixed points of the iteration (15) are maximizers of the log-likelihood function  $l_y(\theta)$ .*

**Proof 3** *Consider a fixed point  $\theta^*$  of iteration (15) for  $\beta_k = \beta = \text{constant}$ . Then,*

$$\theta^* = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \{l_y(\theta) - \beta I(\theta, \theta^*)\}.$$

*As  $l_y(\theta)$  and  $I(\theta, \theta^*)$  are both smooth in  $\theta$ ,  $\theta^*$  must be a stationary point*

$$0 = \nabla l_y(\theta^*) - \beta \nabla_{10} I(\theta^*, \theta^*).$$

*Furthermore, as  $I(\theta, \bar{\theta})$  has a minimum (zero) at  $\theta = \bar{\theta}$ ,  $\nabla_{10} I(\theta^*, \theta^*) = 0$ ,*

$$0 = \nabla l_y(\theta). \tag{17}$$

*Since  $l_y(\theta)$  is strictly concave, we deduce that  $\theta^*$  is a maximizer of  $l_y(\theta)$ .*

**Remark 1** *For nonconvex log-likelihood, equation (17) asserts that fixed points of iteration (15) for constant  $\beta_k$  are stationary points of  $l_y(\theta)$ . Stationary points can be minimizers, maximizers or saddle points.*

The following will be useful.

**Lemma 1** *Let  $k(x; \theta)$  be such that Assumptions 1 are satisfied for  $I(\theta, \bar{\theta})$ . Then, given two bounded sequences  $\{\theta_1^k\}$  and  $\{\theta_2^k\}$ ,  $\lim_{k \rightarrow \infty} I(\theta_1^k, \theta_2^k) = 0$  implies that  $\lim_{k \rightarrow \infty} \|\theta_1^k - \theta_2^k\| = 0$ .*

**Proof 4** *Let  $\mathcal{B}$  be any bounded set containing both sequences  $\{\theta_1^k\}$  and  $\{\theta_2^k\}$ . Let  $\lambda$  denote the minimum*

$$\lambda = \min_{\theta, \bar{\theta} \in \mathcal{B}} \lambda_{\nabla_{10}^2 I(\theta, \bar{\theta})} \tag{18}$$

*Assumption 1.iv implies that  $\lambda > 0$ . Furthermore, invoking Taylor's theorem with remainder,  $I(\theta, \bar{\theta})$  is strictly convex in the sense that for any  $k$*

$$\begin{aligned} I(\theta_1^k, \theta_2^k) &\geq I(\theta_1^k, \theta_1^k) + \nabla I(\theta_1^k, \theta_1^k)^\top (\theta_1^k - \theta_2^k) \\ &\quad + \frac{1}{2} \lambda \|\theta_1^k - \theta_2^k\|^2. \end{aligned}$$

*As  $I(\theta_1^k, \theta_1^k) = 0$  and  $\nabla_{10} I(\theta_1^k, \theta_1^k) = 0$  we obtain*

$$I(\theta_1^k, \theta_2^k) \geq \frac{\lambda}{2} \|\theta_1^k - \theta_2^k\|^2.$$

*The desired result comes from passing to the limit  $k \rightarrow \infty$ .*

Using these results, we easily obtain the following.

**Lemma 2** *Let  $\{g(y; \theta)\}_{\theta \in \mathbb{R}^p}$  and  $\{k(x; \theta)\}_{\theta \in \mathbb{R}^p}$  be such that Assumptions 1 are satisfied. Then  $\{\theta^k\}_{k \in \mathbb{N}}$  is bounded.*

**Proof 5** *Due to Proposition 2, the sequence  $\{l_y(\theta^k)\}$  is monotone increasing. Therefore, assumption 1.ii implies that  $\{\theta^k\}$  is bounded.*

In the following lemma, we prove a result which is often called asymptotic regularity [20].

**Lemma 3** *Let the densities  $\{g(y; \theta)\}_{\theta \in \mathbb{R}^p}$  and  $\{k(x; \theta)\}_{\theta \in \mathbb{R}^p}$  be such that  $l_y(\theta)$  and  $I(\theta, \bar{\theta})$  satisfy Assumptions 1. Assume in addition that  $\{\beta_k\}_{k \in \mathbb{N}}$  converges to  $\beta^* > 0$  for some  $\beta^*$ . Then,*

$$\lim_{k \rightarrow \infty} \|\theta^{k+1} - \theta^k\| = 0. \tag{19}$$

**Proof 6** By Assumption 1.iii and by Proposition 2  $\{l_y(\theta^k)\}_{k \in \mathbb{N}}$  is bounded and monotone. Since, by Lemma 2,  $\{\theta^k\}_{k \in \mathbb{N}}$  is a bounded sequence  $\{l_y(\theta^k)\}_{k \in \mathbb{N}}$  converges. Therefore,  $\lim_{k \rightarrow \infty} \{l_y(\theta^{k+1}) - l_y(\theta^k)\} = 0$  which, from (16), implies

$$\sum_{k \in \mathbb{N}} \beta_k I(\theta^{k+1}, \theta^k) = 0. \quad (20)$$

As the summand is nonnegative, we deduce that  $\beta_k I(\theta^{k+1}, \theta^k)$  vanishes when  $k$  tends to infinity. Since  $\{\beta_k\}_{k \in \mathbb{N}}$  is bounded below by  $\beta^* > 0$ ,

$$\lim_{k \rightarrow \infty} I(\theta^{k+1}, \theta^k) = 0. \quad (21)$$

Therefore, Lemma 1 establishes the desired result.

We can now give a global convergence theorem.

**Theorem 1** For any positive convergent sequence of relaxation parameters  $\{\beta_k\}_{k \in \mathbb{N}}$ , the sequence  $\{\theta^k\}_{k \in \mathbb{N}}$  converges to the solution of the ML estimation problem (1).

**Proof 7** Since  $\{\theta^k\}_{k \in \mathbb{N}}$  is bounded, one can extract a convergent subsequence  $\{\theta^{\sigma(k)}\}_{k \in \mathbb{N}}$  with limit  $\theta^*$ . The defining iteration (15) implies that

$$\nabla l_y(\theta^{\sigma(k)+1}) - \beta_{\sigma(k)} \nabla_{10} I(\theta^{\sigma(k)+1}, \theta^{\sigma(k)}) = 0.$$

We now prove that  $\theta^*$  is a stationary point of  $l_y(\theta)$ . Assume first that  $\{\beta_k\}_{k \in \mathbb{N}}$  converges to zero, i.e.  $\beta^* = 0$ . Due to Assumptions 1.i,  $\nabla l_y(\theta)$  is continuous in  $\theta$ . Hence, since  $\nabla_{10} I(\theta, \bar{\theta})$  is bounded on bounded subsets, (7) implies

$$\nabla l_y(\theta^*) = 0.$$

Next, assume that  $\beta^* > 0$ . In this case, Lemma 3 proves that

$$\lim_{k \rightarrow \infty} \|\theta^{k+1} - \theta^k\| = 0.$$

Therefore,  $\{\theta^{\sigma(k)+1}\}_{k \in \mathbb{N}}$  also tends to  $\theta^*$ . Since  $\nabla_{10} I(\theta, \bar{\theta})$  is continuous in  $(\theta, \bar{\theta})$  equation (7) gives at infinity

$$\nabla l_y(\theta^*) - \beta^* \nabla_{10} I(\theta^*, \theta^*) = 0.$$

Now, it is straightforward that  $\nabla_{10} I(\theta^*, \theta^*) = 0$ . Then,

$$\nabla l_y(\theta^*) = 0. \quad (22)$$

The proof is concluded as follows. As, by Assumption 1.iii,  $l_y(\theta)$  is concave,  $\theta^*$  is a maximizer of  $l_y(\theta)$  so that  $\theta^*$  solves the Maximum Likelihood estimation problem (1). Furthermore, as positive definiteness of  $\nabla^2 l_y$  implies that  $l_y(\theta)$  is in fact strictly concave, this maximizer is unique. Hence,  $\{\theta^k\}$  has only one accumulation point and  $\{\theta^k\}$  converges to  $\theta^*$  which ends the proof.

We now establish the main result concerning speed of convergence.

**Theorem 2** Assume that the sequence of relaxation parameters  $\{\beta_k\}_{k \in \mathbb{N}}$  vanishes. Then,  $\{\theta^k\}$  converges superlinearly to the solution of the ML estimation problem (1).

**Proof 8** Due to Theorem 1, the sequence  $\{\theta^k\}$  converges to the unique maximizer  $\theta_{ML}$  of  $l_y(\theta)$ . Assumption 1.i implies that the gradient mapping  $\nabla_{\theta}(l_y(\theta) - \beta_k I(\theta, \theta_{ML}))$  is continuously differentiable. Hence, we have the following Taylor expansion about  $\theta_{ML}$ .

$$\begin{aligned} \nabla l_y(\theta) - \beta_k \nabla_{10} I(\theta, \theta_{ML}) &= \nabla l_y(\theta_{ML}) \\ &\quad - \beta_k \nabla_{10} I(\theta_{ML}, \theta_{ML}) \\ &\quad + \nabla^2 l_y(\theta_{ML})(\theta - \theta_{ML}) \\ &\quad - \beta_k \nabla_{10}^2 I(\theta_{ML}, \theta_{ML})(\theta - \theta_{ML}) \\ &\quad + R(\theta - \theta_{ML}), \end{aligned} \quad (23)$$

where the remainder satisfies

$$\lim_{\theta \rightarrow \theta_{ML}} \frac{\|R(\theta - \theta_{ML})\|}{\|\theta - \theta_{ML}\|} = 0.$$

Since  $\theta_{ML}$  maximizes  $l_y(\theta)$ ,  $\nabla l_y(\theta_{ML}) = 0$ . Furthermore,  $\nabla_{10} I(\theta_{ML}, \theta_{ML}) = 0$ . Hence, (23) can be simplified to

$$\begin{aligned} \nabla l_y(\theta) - \beta_k \nabla_{10} I(\theta, \theta_{ML}) &= \nabla^2 l_y(\theta_{ML})(\theta - \theta_{ML}) \\ &\quad - \beta_k \nabla_{10}^2 I(\theta_{ML}, \theta_{ML})(\theta - \theta_{ML}) + R(\theta - \theta_{ML}). \end{aligned} \quad (24)$$

From the defining relation (15) the iterate  $\theta^{k+1}$  satisfies

$$\nabla l_y(\theta^{k+1}) - \beta_k \nabla_{10} I(\theta^{k+1}, \theta^k) = 0. \quad (25)$$

So, taking  $\theta = \theta^{k+1}$  in (24) and using (25), we obtain

$$\begin{aligned} &\beta_k (\nabla_{10} I(\theta^{k+1}, \theta^k) - \nabla_{10} I(\theta^{k+1}, \theta_{ML})) = \\ &+ \nabla^2 l_y(\theta_{ML})(\theta^{k+1} - \theta_{ML}) - \beta_k \nabla_{10}^2 I(\theta_{ML}, \theta_{ML})(\theta^{k+1} - \theta_{ML}) \\ &+ R(\theta^{k+1} - \theta_{ML}). \end{aligned}$$

Thus,

$$\begin{aligned} &\|\beta_k (\nabla_{10} I(\theta^{k+1}, \theta^k) - \nabla_{10} I(\theta^{k+1}, \theta_{ML})) - R(\theta^{k+1} - \theta_{ML})\| = \\ &\|\nabla^2 l_y(\theta_{ML})(\theta^{k+1} - \theta_{ML}) - \beta_k \nabla_{10}^2 I(\theta_{ML}, \theta_{ML})(\theta^{k+1} - \theta_{ML})\|. \end{aligned} \quad (26)$$

On the other hand, one deduces from Assumptions 1 (i) that  $\nabla_{10} I(\theta, \bar{\theta})$  is locally Lipschitz in the variables  $\theta$  and  $\bar{\theta}$ . Then, since,  $\{\theta^k\}$  is bounded, there exists a bounded set  $\mathcal{B}$  containing  $\{\theta^k\}$  and a finite constant  $L$  such that for all  $\theta, \theta', \bar{\theta}$  and  $\bar{\theta}'$  in  $\mathcal{B}$ ,

$$\|\nabla_{10} I(\theta, \bar{\theta}) - \nabla_{10} I(\theta', \bar{\theta}')\| \leq L(\|\theta - \theta'\|^2 + \|\bar{\theta} - \bar{\theta}'\|^2)^{\frac{1}{2}}.$$

Using the triangle inequality and this last result, (26) asserts that for any  $\theta \in \mathcal{B}$

$$\begin{aligned} &\beta_k L \|\theta^{k+1} - \theta_{ML}\| + \|R(\theta^{k+1} - \theta_{ML})\| \geq \|(\nabla^2 l_y(\theta_{ML}) \\ &\quad - \beta_k \nabla_{10}^2 I(\theta_{ML}, \theta_{ML}))(\theta^{k+1} - \theta_{ML})\|. \end{aligned} \quad (27)$$

Now, consider again the bounded set  $\mathcal{B}$  containing  $\{\theta^k\}$ . Let  $\lambda_{l_y}$  and  $\lambda_I$  denote the minima

$$\lambda_{l_y} = \min_{\theta \in \mathcal{B}} \{-\lambda_{\nabla^2 l_y(\theta)}\}$$

$$\lambda_I = \min_{\theta, \bar{\theta} \in \mathcal{B}} \{\lambda_{\nabla_{10}^2 I(\theta, \bar{\theta})}\}.$$

Since for any symmetric matrix  $H$ ,  $x^T H x / \|x\|^2$  is lower bounded by the minimum eigenvalue of  $H$ , we have immediately that

$$\begin{aligned} &\|(-\nabla^2 l_y(\theta_{ML}) + \beta_k \nabla_{10}^2 I(\theta_{ML}, \theta_{ML}))(\theta^{k+1} - \theta_{ML})\|^2 \\ &\geq (\lambda_{l_y} + \beta_k \lambda_I)^2 \|\theta^{k+1} - \theta_{ML}\|^2. \end{aligned} \quad (28)$$

By Assumptions 1.iii and 1.iv,  $\lambda_{l_y} + \beta_k \lambda_I > 0$  and, after substitution of (28) into (27), we obtain

$$\begin{aligned} &\beta_k L \|\theta^k - \theta_{ML}\| + \|R(\theta^{k+1} - \theta_{ML})\| \geq \\ &\quad (\lambda_{l_y} + \beta_k \lambda_I) \|\theta^{k+1} - \theta_{ML}\|, \end{aligned} \quad (29)$$

for all  $\theta \in \mathcal{B}$ . Therefore, collecting terms in (29)

$$\beta_k L \geq \left( \lambda_{l_y} + \beta_k \lambda_I - \frac{\|R(\theta^{k+1} - \theta_{ML})\|}{\|\theta^{k+1} - \theta_{ML}\|} \right) \frac{\|\theta^{k+1} - \theta_{ML}\|}{\|\theta^k - \theta_{ML}\|}. \quad (30)$$

Now, recall that  $\{\theta^k\}$  is convergent. Thus,  $\lim_{k \rightarrow \infty} \|\theta^k - \theta_{ML}\| = 0$  and subsequently,  $\lim_{k \rightarrow \infty} \frac{\|R(\theta^{k+1} - \theta_{ML})\|}{\|\theta^{k+1} - \theta_{ML}\|} = 0$  due to the definition of the remainder  $R$ . Since convergence of  $\{\theta^k\}$  also implies that  $\lim_{k \rightarrow \infty} \|\theta^k - \theta^{k-1}\| = 0$ , invoking the assumptions of Lemma 3 yields:  $\lim_{k \rightarrow \infty} \beta_k = \lim_{k \rightarrow \infty} \phi(\|\theta^k - \theta^{k-1}\|) = 0$ . Therefore, as  $\lambda_{l_y} > 0$ , equation (30) gives

$$\lim_{k \rightarrow \infty} \frac{\|\theta^{k+1} - \theta_{ML}\|}{\|\theta^k - \theta_{ML}\|} = 0,$$

and the proof of superlinear convergence is completed.

**Remark 2** The convergence Theorems 1 and 2 make use of concavity of  $l_y(\theta)$  and convexity of  $I(\theta, \bar{\theta})$  via Assumptions 1.iii and 1.iv. However, for smooth non-convex functions an analogous local superlinear convergence result can be established under stronger assumptions similar to those in [3].

**Remark 3** The convergence Theorems 1 and 2 also clearly apply to a class of objective functions which is not restricted to Likelihood functions. For instance, the analysis directly adapts to Penalized Maximum Likelihood problems.

## 5 Second order Approximations and Trust Region techniques

The main drawback of the Kullback proximal point algorithm (15) is that for  $\beta_k \neq 1$ , for which it reduces to EM, recursion (15) may be difficult to implement. In this section, we discuss an easily implementable version of the KPP algorithm using second order approximations which preserve monotonicity. The second order scheme is related to the well-known Trust Region technique for optimization introduced by Moré [21].

### 5.1 Approximate models

In order to obtain computable iterations, the following second order approximations of  $l_y(\theta)$  and  $I(\theta, \theta^k)$  can be introduced

$$\begin{aligned} \hat{l}_y(\theta) &= l_y(\theta^k) + \nabla l_y(\theta^k)^\top (\theta - \theta^k) + \\ &\quad \frac{1}{2} (\theta - \theta^k)^\top \nabla^2 l_y(\theta) (\theta - \theta^k). \end{aligned}$$

and

$$\hat{I}(\theta, \theta^k) = \frac{1}{2} (\theta - \theta^k)^\top \nabla_{10}^2 I(\theta^k, \theta^k) (\theta - \theta^k).$$

In the following, we adopt the notations

$$\begin{aligned} g_k &= \nabla l_y(\theta^k) \\ H_k &= \nabla^2 l_y(\theta^k) \end{aligned}$$

and

$$I_k = \nabla_{10}^2 I(\theta^k, \theta^k).$$

The approximate KPP algorithm is defined as

$$\begin{aligned} \theta^{k+1} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \{ & l_y(\theta^k) + g_k(\theta - \theta^k) \\ & + \frac{1}{2}(\theta - \theta^k)^\top H_k(\theta - \theta^k) \\ & - \frac{\beta_k}{2}(\theta - \theta^k)^\top I_k(\theta - \theta^k) \} \end{aligned} \quad (31)$$

At this point it is important to make several comments. Notice first that for  $\beta_k = 0$ ,  $k = 1, 2, \dots$ , the approximate step (31) is equivalent to a Newton step. It is well known that Newton's method, also known as Fisher scoring in statistics, has superlinear asymptotic convergence rate but may diverge if not properly initialized. Therefore, at least for small values of the relaxation parameter  $\beta_k$ , the approximate PPA algorithm may fail to converge for the same reasons as for Newton's method. On the other hand, for  $\beta_k > 0$  the term  $-\frac{\beta_k}{2}(\theta - \theta^k)^\top I_k(\theta - \theta^k)$  penalizes the distance of the next iterate  $\theta^{k+1}$  to the current iterate  $\theta^k$ . Hence, we can interpret this term as a regularization which stabilizes the possibly divergent Newton algorithm without sacrificing its superlinear asymptotic convergence rate. By appropriate choice of  $\{\beta_k\}$  the iterate  $\theta^{k+1}$  can be forced to remain in a region around  $\theta^k$  over which the quadratic model  $\hat{l}_y(\theta)$  is accurate. This idea is very popular in optimization and is better known as the "Trust Region technique" [21][22].

**Remark 4** *In many cases a quadratic approximation of a single one of the two terms  $l_y(\theta)$  or  $I(\theta, \theta^k)$  is sufficient to obtain a closed form maximization (14). Naturally, when feasible, such a reduced approximation is preferable to the approximation of both terms discussed above. For concreteness, in the sequel, although our results hold for the reduced approximation also, we prove convergence for the proximal point algorithm implemented with the full two-term approximation only.*

## 5.2 Trust Region Techniques

The Trust Region strategy proceeds as follows. The model  $\hat{l}_y(\theta)$  is maximized in a ball  $B(\theta^k, \delta) = \{\|\theta - \theta^k\|_{I_k} \leq \delta\}$  centered at  $\theta^k$  where  $\delta$  is a proximity control parameter which may depend on  $k$ , and where  $\|a\|_{I_k} = a^\top I_k a$  is a norm; well defined due to positive definiteness of  $I_k$  (Assumption 1.iv). Given an iterate  $\theta^k$  consider a candidate  $\theta^\delta$  for  $\theta^{k+1}$  defined as the solution to the constrained optimization problem

$$\theta^\delta = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \hat{l}_y(\theta)$$

subject to

$$\|\theta - \theta^k\|_{I_k} \leq \delta. \quad (32)$$

By duality theory of constrained optimization [23], and the fact that  $\hat{l}_y(\theta)$  is strictly concave, this problem is equivalent to the unconstrained optimization

$$\theta^\delta(\beta) = \operatorname{argmin}_{\theta \in \mathbb{R}^p} L(\theta, \beta). \quad (33)$$

where

$$L(\theta, \beta) = -\hat{l}_y(\theta) + \frac{\beta}{2}(\|\theta - \theta^k\|_k^2 - \delta^2).$$

and  $\beta$  is a Lagrange multiplier selected to meet the constraint (32) with equality:  $\|\theta^\delta(\beta) - \theta\|_{I_k} = \delta$ . We conclude that the Trust Region candidate  $\theta^\delta$  is identical to the approximate KPP iterate (31) with regularization parameter  $\beta$  chosen according to constraint (32). This relation also provides a rational rule for computing the relaxation parameter  $\beta$ .

### 5.3 Implementation

The parameter  $\delta$  is said to be safe if  $\theta^\delta$  produces an acceptable increase in the original objective  $l_y$ . An iteration of the Trust Region method consists of two principal steps

*Rule 1.* Determine whether  $\delta$  is safe or not. If  $\delta$  is safe, set  $\delta_k = \delta$  and take an approximate Kullback proximal step  $\theta^{k+1} = \theta^\delta$ . Otherwise, take a *null step*  $\theta^{k+1} = \theta^k$ .

*Rule 2.* Update  $\delta$  depending on the result of *Rule 1*.

Rule 1 can be implemented by comparing the increase in the original log-likelihood  $l_y$  to a fraction  $m$  of the expected increase predicted by the approximate model  $\hat{l}_y(\theta)$ . Specifically, the Trust Region parameter  $\delta$  is accepted if

$$l_y(\theta^\delta) - l_y(\theta^k) \geq m(\hat{l}_y(\theta^\delta) - \hat{l}_y(\theta^k)). \quad (34)$$

Rule 2 can be implemented as follows. If  $\delta$  was accepted by Rule 1,  $\delta$  is increased at the next iteration in order to extend the region of validity of the model  $\hat{l}_y(\theta)$ . If  $\delta$  was rejected, the region must be tightened and  $\delta$  is decreased at the next iteration.

The Trust Region strategy implemented here is essentially the same as that proposed by Moré [21].

**Algorithm 1** *Step 0. (Initialization)* Set  $\theta^0 \in \mathbb{R}^p$ ,  $\delta_0 > 0$  and the “curve search” parameters  $m$ ,  $m'$  with  $0 < m < m' < 1$ .

*Step 1. Solve*

$$\theta^{\delta_k} = \operatorname{argmax}_{\theta \in \mathbb{R}^p} \hat{l}_y(\theta)$$

*subject to*

$$\|\theta - \theta^k\| \leq \delta_k.$$

*Step 2.* If  $l_y(\theta^{\delta_k}) - l_y(\theta^k) \geq m(\hat{l}_y(\theta^{\delta_k}) - \hat{l}_y(\theta^k))$  then set  $\theta^{k+1} = \theta^{\delta_k}$ . Otherwise, set  $\theta^{k+1} = \theta^k$ .

*Step 3.* Set  $k = k + 1$ . Update the model  $\hat{l}_y^k(\theta)$ . Update  $\delta_k$  using Procedure 1.

*Step 4.* Go to Step 1.

The procedure for updating  $\delta_k$  is given below.

**Procedure 1** *Step 0. (Initialization)* Set  $\gamma_1$  and  $\gamma_2$  such that  $\gamma_1 < 1 < \gamma_2$ .

*Step 1.* If  $l_y(\theta^{\delta_k}) - l_y(\theta^k) \leq m(\hat{l}_y(\theta^{\delta_k}) - \hat{l}_y(\theta^k))$  then take  $\delta_{k+1} \in (0, \gamma_1 \delta_k)$ .

*Step 2.* If  $l_y(\theta^{\delta_k}) - l_y(\theta^k) \leq m'(\hat{l}_y(\theta^{\delta_k}) - \hat{l}_y(\theta^k))$  then take  $\delta_{k+1} \in (\gamma_1 \delta_k, \delta_k)$ .

*Step 3.* If  $l_y(\theta^{\delta_k}) - l_y(\theta^k) \geq m'(\hat{l}_y(\theta^{\delta_k}) - \hat{l}_y(\theta^k))$  then take  $\delta_{k+1} \in (\delta_k, \gamma_2 \delta_k)$ .

This algorithm satisfies the following convergence theorem

**Theorem 3** Let  $g(y; \theta)$  and  $k(x; \theta)$  be such that Assumptions 1 are satisfied. Then,  $\{\theta^k\}$  generated by Algorithm 1 converges to the maximizer  $\theta_{ML}$  of the log-likelihood  $l_y(\theta)$  and satisfies the monotone likelihood property  $l_y(\theta^{k+1}) \geq l_y(\theta^k)$ . If in addition, the sequence of Lagrange multipliers  $\{\beta_k\}$  tends towards zero,  $\{\theta^k\}$  converges superlinearly.

**Remark 5** The proof of Theorem 3 is omitted since it is standard in the analysis of Trust Region methods. Superlinear convergence as  $\{\beta_k\}$  vanishes comes from the Dennis and More criterion [22, Theorem 3.11].

**Remark 6** The Trust Region framework can also be applied to nonconvex contexts. In the case where  $I_k$  remains positive definite, global convergence to a local maximizer of  $l_y(\theta)$  can be obtained under Assumptions 1.i and 1.ii only following the proof technique of [21].

## 6 Application to Poisson data

In this section, we illustrate the application of Algorithm 1 for a maximum likelihood estimation problem in a Poisson inverse problem arising in emission computed tomography (ECT).

## 6.1 The Poisson Inverse Problem

The objective is to estimate the intensity vector  $\theta = [\theta_1, \dots, \theta_p]^T$  governing the number of gamma-ray emissions  $N = [N_1, \dots, N_p]^T$  over an imaging volume of  $p$  pixels. The estimate of  $\theta$  must be based on a vector of  $m$  observed projections of  $N$  denoted  $Y = [Y_1, \dots, Y_m]^T$ . The components of  $N_i$  of  $N$  are independent Poisson distributed with rate parameters  $\theta_i$ , and the components  $Y_j$  of  $Y$  are independent Poisson distributed with rate parameters  $\sum_{i=1}^p P_{ji}\theta_i$ , where  $P_{ji}$  is the transition probability corresponding to emissions from pixel  $i$  being detected at detector module  $j$ . The standard choice of complete data  $X$ , introduced by Shepp and Vardi [24], is the set  $\{N_{ji}\}_{1 \leq j \leq m, 1 \leq i \leq p}$ , where  $N_{ji}$  denotes the number of emissions in pixel  $i$  which are detected at detector  $j$ . The corresponding many-to-one mapping  $h(X) = Y$  in the EM algorithm is

$$Y_j = \sum_{i=1}^p N_{ji}, \quad 1 \leq j \leq m. \quad (35)$$

It is also well known that the likelihood function is given by

$$\log g(y; \theta) = \sum_{j=1}^m \left( \sum_{i=1}^p P_{ji}\theta_i \right) - y_j \log \left( \sum_{i=1}^p P_{ji}\theta_i \right) + \log y_j! \quad (36)$$

and that the expectation step of the EM algorithm is (see [25])

$$Q(\theta, \bar{\theta}) = \mathbb{E}[\log f(x; \theta) \mid y; \bar{\theta}] = \quad (37)$$

$$\sum_{j=1}^m \sum_{i=1}^p \left( \frac{y_j P_{ji} \bar{\theta}_i}{\sum_{i=1}^p P_{ji} \bar{\theta}_i} \log(P_{ji}\theta_i) - P_{ji}\theta_i \right).$$

## 6.2 Simulation results

For illustration of the convergence properties of the algorithms studied here, we performed simulations on a simple one dimensional deblurring example with Poisson noise model. This example easily generalizes to more general 2 and 3 dimensional Poisson deblurring, tomographic reconstruction, and other imaging applications. The true source  $\theta$  is a two level phantom shown in Figure 1. The blurring kernel is a Gaussian function yielding the noiseless blurred phantom shown in Figure 2. In Figure 5 are shown the results of deblurring in the absence of noise via optimization of the  $Q$  function 37. Our simulation results show in Figure 3 that, the Trust Region implementation of our Kullback proximal algorithm enjoys much faster convergence towards the optimum than EM after only a few iterations, for  $(\beta_k)_{k \in \mathbb{N}}$  being the sequence of Lagrange parameters associated to the sequence of trust region radii  $(\delta_k)_{k \in \mathbb{N}}$  updated as in Procedure 1. Figure 4 validates the theoretical superlinear convergence of the Trust Region iterates as contrasted with the linear convergence rate of the EM iterates. Figure 5 shows the reconstructed signals and demonstrates that the Trust Region technique achieves better reconstruction of the original phantom. Finally, Figure 6 shows the iterates for the reconstructed phantom, plotted as a function of iteration on the horizontal axis and as a function of grey level on the vertical axis. Observe that the proximal point EM achieves more rapid separation of the two components in the phantom than does standard EM.

## 7 Conclusions

The main contributions of this paper are the following. First, we introduced a very general class of iterative methods for ML estimation based on Kullback regularization and the proximal point strategy. Next, we proved that the EM algorithm belongs to the proposed class, thus providing a new and insightful interpretation of the EM approach for ML estimation. The proximal approach developed here naturally adapts to penalized Maximum Likelihood problems. Finally, we showed that Kullback proximal point methods enjoy global convergence and even superlinear convergence for vanishing sequences of relaxation parameters. Implementation issues were also discussed and we provided second order schemes for the case



where the Maximization step is hard to obtain in closed form. We addressed Trust Region methodologies for the updating of the relaxation parameter. Computational experiments provided evidences for the good behavior of the approximate second order scheme in practice and showed superlinear behavior of the iterates on a simple example. Several extensions of the proposed approach are currently under investigation. In particular a Proximal generalization of the SAGE algorithm of Fessler and Hero [26] is proposed in [27], in the case of mixture densities estimation.

## References

- [1] A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum likelihood for incomplete data via the EM algorithm (with discussion),” *Journal of the Royal Statistical Society*, vol. 39, pp. 1–38, 1977.
- [2] C. F. Wu, “On the convergence of the em algorithm,” *Ann. Statis.*, vol. 11, no. 1, pp. 95–103, 1983.
- [3] A. O. Hero and J. A. Fessler, “Convergence in norm for em-type algorithms,” *Statistica Sinica*, vol. 5, no. 1, pp. 41–54, 1995.
- [4] A. K. Louis, “Nonuniqueness problems in computerized tomography,” *Z. Angew. Math Mech.*, vol. 62, pp. 290–292, 1982.
- [5] R. M. Lewitt and Muehlehner, “Accelerated iterative reconstruction for positron emission tomography based on the em algorithm for maximum likelihood estimation,” *IEEE Tr. Med. Im.*, vol. 5, no. 1, pp. 16–22, 1986.
- [6] L. Kaufman, “Implementing and accelerating the em algorithm for positron emission tomography,” *IEEE Tr. Med. Im.*, vol. 6, no. 1, pp. 37–51, 1987.
- [7] M. Jamshidian and R. I. Jennrich, “Conjugate gradient acceleration of the em algorithm,” *J. Am. Stat. Ass.*, vol. 88, no. 421, pp. 221–228, 1993.
- [8] I. Meilijson, “A fast improvement to the em algorithm in its own terms,” *J. Roy. Stat. Soc. Ser. B*, vol. 51, no. 1, pp. 127–138, 1989.
- [9] C. Bouman and K. Sauer, “A unified approach to statistical tomography using coordinate descent optimization,” *Proc. 27th Conf. Info. Sci. Sys., John Hopkins*, pp. 611–616, 1993.
- [10] Lange, Kenneth, “A quasi-Newton acceleration of the EM algorithm.,” *Stat. Sin. 5, No.1, 1-18 (1995)*.
- [11] H M Hudson and R S Larkin, “Accelerated image reconstruction using ordered subsets of projection data,” *IEEE Tr. Med. Im.*, vol. 13, no. 4, pp. 601, Dec. 1994.
- [12] Lavielle, Marc, “A stochastic algorithm for parametric and non-parametric estimation in the case of incomplete data.,” *Signal Process. 42, No.1, 3-17 (1995). [ISSN 0165-1684 ]*.
- [13] D. Lansky and G. Casella, “Improving the EM algorithm,” in *Computing and Statistics: Proc. Symp. on the Interface*, C. Page and R. LePage, Eds., pp. 420–424. Springer-Verlag, 1990.
- [14] R. T. Rockafellar, “Monotone operators and the proximal point algorithm,” *SIAM Journal on Control and Optimization*, vol. 14, pp. 877–898, 1976.
- [15] B. Martinet, “Régularisation d’inéquation variationnelles par approximations successives,” *Revue Francaise d’Informatique et de Recherche Operationnelle*, vol. 3, pp. 154–179, 1970.
- [16] J. B. Hiriart Urruty and C. Lemaréchal, *Convex Analysis and Minimization Algorithms I-II*, Springer Verlag, 1993, Grundlehren der mathematischen Wissenschaften 306.
- [17] M. Teboulle, “Entropic proximal mappings with application to nonlinear programming,” *Mathematics of Operations Research*, vol. 17, pp. 670–690, 1992.

- [18] I. A. Ibragimov and R. Z. Khas'minskij, *Asymptotic theory of estimation*, Teoriya Veroyatnostej i Matematicheskaya Statistika. Moskva:.
- [19] M. Teboulle, "Convergence of proximal-like algorithms," *SIAM Journal on Optimization*, vol. 7, pp. 1069–1083, 1997.
- [20] H. H. Bauschke and J. M. Borwein, "On projection algorithms for solving convex feasibility problems," *SIAM Review*, vol. 38, no. 3, pp. 367–426, 1996.
- [21] J. J. Moré, "Recent developments in algorithms and software for trust region methods," *Mathematical Programming: The State of the Art, Bonn, Springer Verlag*, pp. 258–287, 1983.
- [22] J. F. Bonnans, J.-Ch. Gilbert, C. Lemaréchal, and C. Sagastizàbal, *Optimization numérique. Aspects théoriques et pratiques*, vol. 27, Springer Verlag, 1997, Series : Mathématiques et Applications.
- [23] J. B. Hiriart Urruty and C. Lemaréchal, *Convex Analysis and Minimization Algorithms I-II*, Springer Verlag, 1993, Grundlehren der mathematischen Wissenschaften 306.
- [24] Vardi, Y. and Shepp, L.A. and Kaufman, L. , "A statistical model for positron emission tomography," *J. Am. Stat. Assoc.* 80, 8-37 (1985).
- [25] Green, Peter J., "On use of EM algorithm for penalized likelihood estimation.," *J. R. Stat. Soc., Ser. B* 52, No.3, 443-452 (1990).
- [26] J. A. Fessler and A. O. Hero, "Space-alternating generalized expectation-maximisation algorithm," *IEEE Trans. Signal Processing*, vol. 42, pp. 2664–2677, 1994.
- [27] G. Celeux, S. Chretien, F. Forbes, and A. Mkhadri, "A component-wise EM algorithm for mixtures," *Research Report 3746 INRIA*, <http://www.inria.fr/RRRT/publications-fra.html>, Aug. 1999.

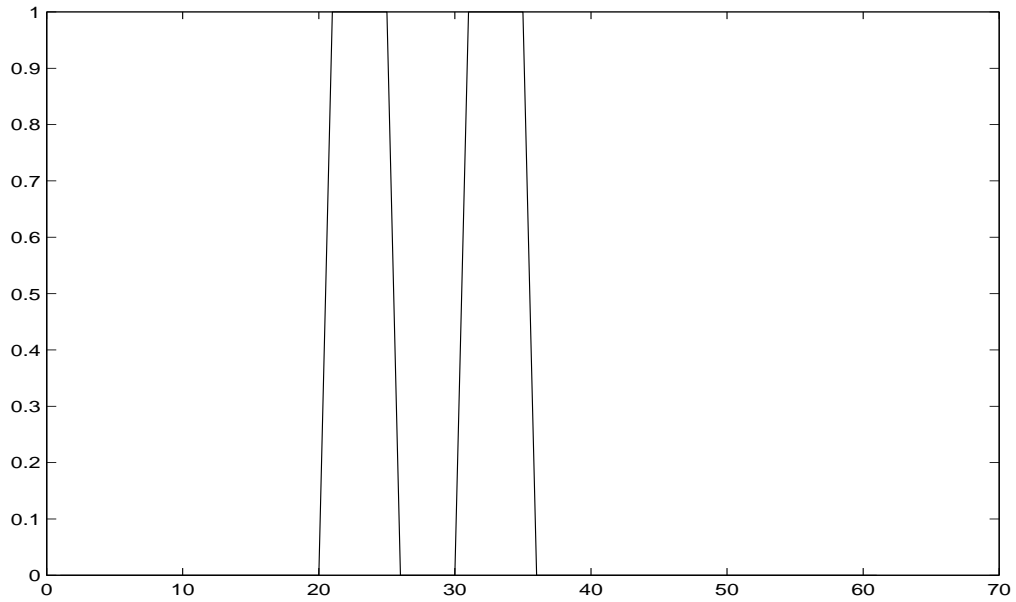


Figure 1: Original phantom

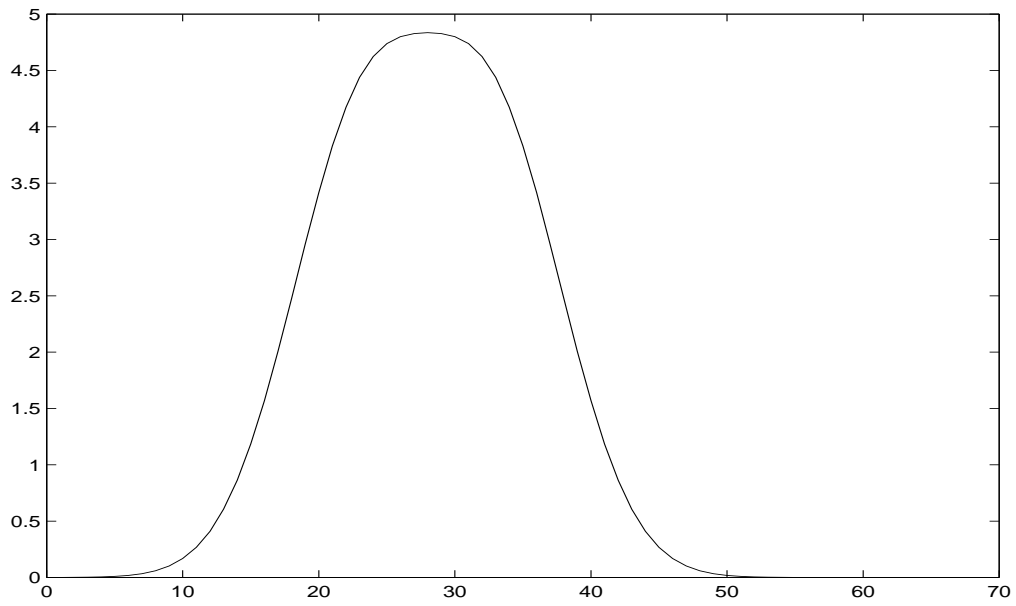


Figure 2: Blurred phantom



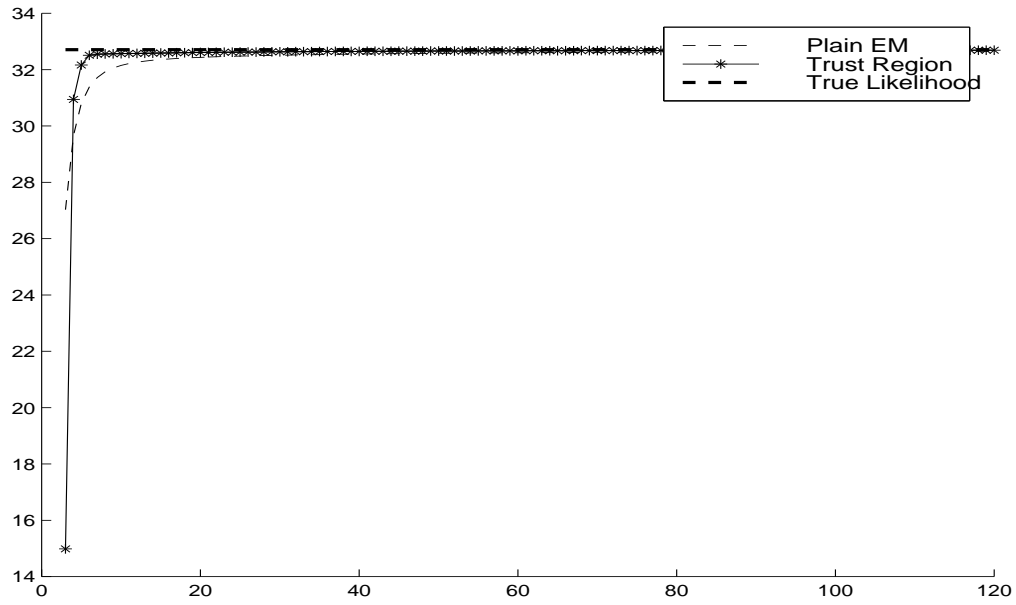
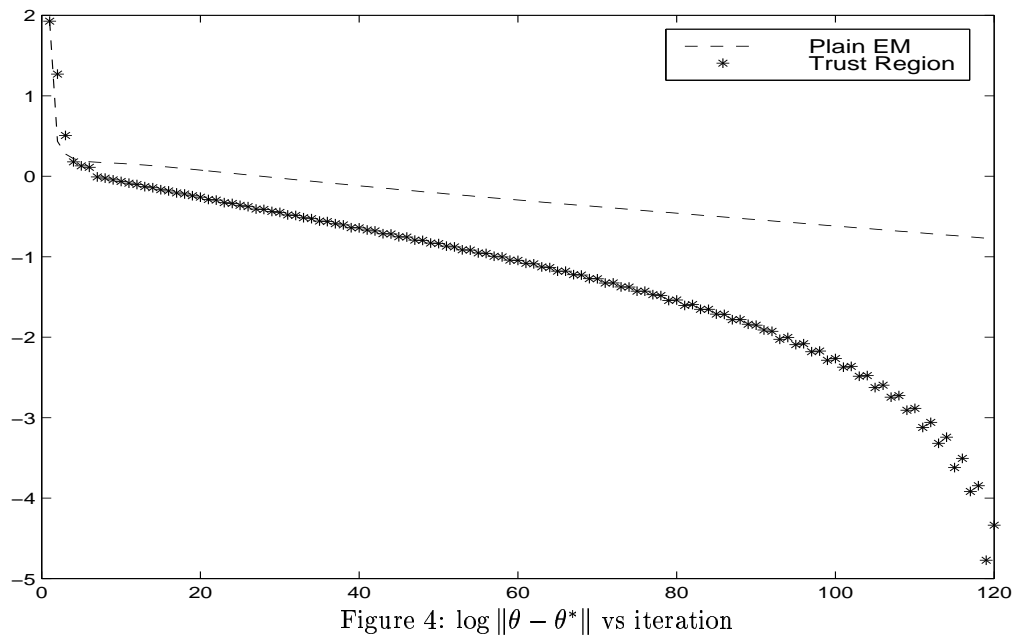


Figure 3: Log-Likelihood vs iteration

Figure 4:  $\log \|\theta - \theta^*\|$  vs iteration

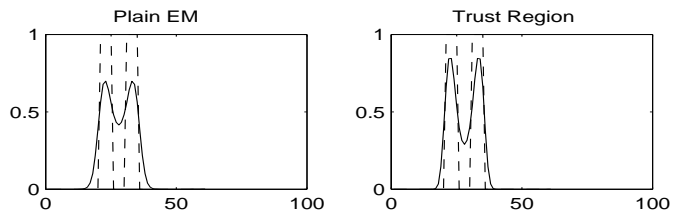


Figure 5: Reconstructed signals after 120 iterations

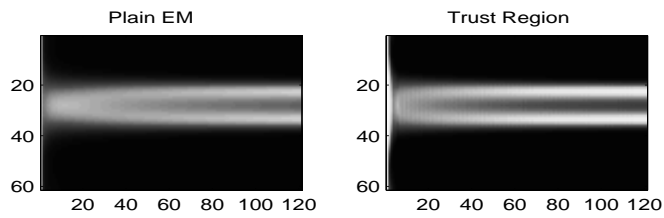


Figure 6: Evolution of the reconstructed image vs iteration

Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,  
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY  
Unité de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex  
Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN  
Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex  
Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

---

Éditeur  
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)  
<http://www.inria.fr>  
ISSN 0249-6399