

Choosing parameters in block-iterative or ordered subset reconstruction algorithms

Charles Byrne (Charles_Byrne@uml.edu),
Department of Mathematical Sciences,
University of Massachusetts Lowell, Lowell, MA 01854

May 5, 2004

Abstract

Viewed abstractly, all the algorithms considered here are designed to provide a nonnegative solution x to the system of linear equations $y = Px$, where y is a vector with positive entries and P a matrix whose entries are nonnegative and with no purely zero columns. The *expectation maximization maximum likelihood* (EMML) method as it occurs in emission tomography and the *simultaneous multiplicative algebraic reconstruction technique* (SMART) are slow to converge on large data sets; accelerating convergence through the use of block-iterative or ordered subset versions of these algorithms is a topic of considerable interest. These block-iterative versions involve relaxation and normalization parameters the correct selection of which may not be obvious to all users. The algorithms are not faster merely by virtue of being block-iterative; the correct choice of the parameters is crucial. Through a detailed discussion of the theoretical foundations of these methods we come to a better understanding of the precise roles these parameters play.

To appear in IEEE Transactions on Image Processing.

1 Introduction

Image reconstruction problems in tomography are often formulated as statistical likelihood maximization problems in which the pixel values of the desired image play the role of parameters. Iterative algorithms based on cross-entropy minimization, such as the *expectation maximization maximum likelihood* (EMML) method [1] and the *simultaneous multiplicative algebraic reconstruction technique* (SMART) [2, 3, 4, 5, 6] can be used to solve such problems. Because the EMML and SMART are slow to converge for the large data sets typical in imaging problems acceleration of the algorithms using blocks of data or ordered subsets has become popular. There are a number of different ways to formulate these block-iterative versions of EMML and SMART, involving the

choice of certain normalization and relaxation parameters. These methods are not faster merely because they are block-iterative; the correct choice of the parameters is crucial. The purpose of this paper is to discuss these different formulations in detail sufficient to reveal the precise roles played by the parameters and to guide the user in choosing them. This is not a survey of the field of iterative algorithms and no attempt has been made to give complete references for each algorithm mentioned here. The reader should consult the editorial [7] for a fuller list of references to the literature.

The notion of *cross-entropy* or the *Kullback-Leibler distance* [8] is central to our discussion. For positive numbers a and b let

$$KL(a, b) = a \log \frac{a}{b} + b - a;$$

also let $KL(a, 0) = +\infty$ and $KL(0, b) = b$. It is easily seen that $KL(a, b) > 0$ unless $a = b$. We extend this Kullback-Leibler distance component-wise to vectors x and z with nonnegative entries:

$$KL(x, z) = \sum_{j=1}^J KL(x_j, z_j). \quad (1.1)$$

The KL distance is also often called Csiszár's I-divergence, denoted $I(a||b)$ [9, 10]. Note that $KL(x, z)$ and $KL(z, x)$ are generally not the same. While the KL distance is not a metric in the usual sense it does have certain properties involving best approximation that are similar to those of the square of the Euclidean metric.

The methods based on cross-entropy, such as the multiplicative version of the *algebraic reconstruction technique* (ART), the MART [11], its simultaneous version, SMART, the expectation maximization maximum likelihood method (EMML) and all block-iterative versions of these algorithms apply to nonnegative systems that we denote by $Px = y$, where y is a vector of positive entries, P is a matrix with entries $P_{ij} \geq 0$ such that for each j the sum $s_j = \sum_{i=1}^I P_{ij}$ is positive and we seek a solution x with nonnegative entries. If no nonnegative x satisfies $y = Px$ we say the system is *inconsistent*.

Simultaneous iterative algorithms employ all of the equations at each step of the iteration; block-iterative methods do not. For the latter methods we assume that the index set $\{i = 1, \dots, I\}$ is the (not necessarily disjoint) union of the N sets or *blocks* B_n , $n = 1, \dots, N$. We shall require that $s_{nj} = \sum_{i \in B_n} P_{ij} > 0$ for each n and each j . Block-iterative methods like ART and MART for which each block consists of precisely one element are called *row-action* or *sequential* methods. We begin our discussion with the SMART and the EMML method.

2 The SMART and the EMMML methods

Both the SMART and the EMMML method provide a solution of $y = Px$ when such exist and (distinct) approximate solutions in the inconsistent case. Both begin with an arbitrary positive vector x^0 . Having found x^k the iterative step for the SMART is

SMART:

$$x_j^{k+1} = x_j^k \exp \left(s_j^{-1} \sum_{i=1}^I P_{ij} \log \frac{y_i}{(Px^k)_i} \right) \quad (2.1)$$

while that for the EMMML method is

EMMML:

$$x_j^{k+1} = x_j^k s_j^{-1} \sum_{i=1}^I P_{ij} \frac{y_i}{(Px^k)_i}. \quad (2.2)$$

The following theorem summarizes what we know of SMART from the references above.

Theorem 2.1 *In the consistent case the SMART converges to the unique nonnegative solution of $y = Px$ for which the distance $\sum_{j=1}^J s_j KL(x_j, x_j^0)$ is minimized. In the inconsistent case it converges to the unique nonnegative minimizer of the distance $KL(Px, y)$ for which $\sum_{j=1}^J s_j KL(x_j, x_j^0)$ is minimized; if P and every matrix derived from P by deleting columns has full rank then there is a unique nonnegative minimizer of $KL(Px, y)$ and at most $I - 1$ of its entries are nonzero.*

The next theorem summarizes what we know of the EMMML method from the references above.

Theorem 2.2 *In the consistent case the EMMML algorithm converges to a nonnegative solution of $y = Px$. In the inconsistent case it converges to a nonnegative minimizer of the distance $KL(y, Px)$; if P and every matrix derived from P by deleting columns has full rank then there is a unique nonnegative minimizer of $KL(y, Px)$ and at most $I - 1$ of its entries are nonzero.*

In the consistent case there may be multiple nonnegative solutions and the one obtained using the EMMML algorithm will depend on the starting vector x^0 ; how it depends on x^0 is an open question. These theorems are special cases of more general results on block-iterative methods that we shall prove later in this paper.

Those who have used the SMART or the EMLL on sizable problems have certainly noticed that they are both slow to converge. An important issue, therefore, is how to accelerate convergence. One popular method is through the use of *block-iterative* (or *ordered subset*) methods.

To illustrate block-iterative methods and to motivate our subsequent discussion we consider now the *ordered subset* EM algorithm (OSEM)[12], which is a popular technique in some areas of medical imaging, as well as an analogous version of SMART, which we shall call here the OSSMART. The OSEM algorithm is now used quite frequently in tomographic image reconstruction, where it is acknowledged to produce usable images significantly faster than EMLL method.

The idea behind the OSEM (OSSMART) is simple: the iteration looks very much like the EMLL (SMART), but at each step of the iteration the summations are taken only over the current block. The blocks are processed cyclically.

The OSEM iteration is the following: for $k = 0, 1, \dots$ and $n = k(\text{mod } N) + 1$, having found x^k let

OSEM:

$$x_j^{k+1} = x_j^k s_{nj}^{-1} \sum_{i \in B_n} P_{ij} \frac{y_i}{(Px^k)_i}. \quad (2.3)$$

The OSSMART has the following iterative step:

OSSMART:

$$x_j^{k+1} = x_j^k \exp \left(s_{nj}^{-1} \sum_{i \in B_n} P_{ij} \log \frac{y_i}{(Px^k)_i} \right). \quad (2.4)$$

In general we do not expect block-iterative algorithms to converge in the inconsistent case, but to exhibit *subsequential convergence* to a *limit cycle*, as we shall discuss later. We do, however, want them to converge to a solution in the consistent case; the OSEM and OSSMART do this when the matrix P and the set of blocks $\{B_n, n = 1, \dots, N\}$ satisfy the condition known as *subset balance*, which means that the sums s_{nj} depend only on j and not on n , but not generally. While subset balance may be approximately valid in some special cases it is overly restrictive, eliminating, for example, almost every set of blocks whose cardinalities are not all the same. When the OSEM does well in practice in medical imaging it is probably because the N is not large and only a few iterations are carried out.

The experience with the OSEM was encouraging, however, and strongly suggested that an equally fast, but mathematically correct, block-iterative version of EMLL

could be found; this is the *rescaled block-iterative* EMMML (RBI-EMML)[13]. Both RBI-EMML and an analogous corrected version of OSSMART, the RBI-SMART, provide fast convergence to a solution in the consistent case, for any choice of blocks.

Both the EMMML and SMART are related to likelihood maximization. Minimizing the function $KL(y, Px)$ is equivalent to maximizing the likelihood when the y_i are taken to be measurements of independent Poisson random variables having means $(Px)_i$. The entries of x are the parameters to be determined. This situation arises in emission tomography. So the EMMML is a likelihood maximizer, as its name suggests.

The connection between SMART and likelihood maximization is a bit more convoluted. Suppose that $s_j = 1$ for each j . To minimize $KL(x, x^0)$ subject to $y = Px$ we form the Lagrangian

$$KL(x, x^0) + \sum_{i=1}^I \lambda_i (y_i - (Px)_i), \quad (2.5)$$

and set to zero the partial derivatives with respect to the entries of x . From this we see that the solution necessarily has the form

$$x_j = x_j^0 \exp\left(\sum_{i=1}^I P_{ij} \lambda_i\right) \quad (2.6)$$

for some vector λ with entries λ_i . This *log linear* form also arises in transmission tomography, where it is natural to assume that $s_j = 1$ for each j and $\lambda_i \leq 0$ for each i . We have the following lemma from [2] that helps to connect the SMART algorithm with the transmission tomography problem:

Lemma 2.1 *Minimizing $KL(d, x)$ over x as in (2.6) is equivalent to minimizing $KL(x, x^0)$, subject to $Px = Pd$.*

With $x_+ = \sum_{j=1}^J x_j > 0$ the vector p with entries $p_j = x_j/x_+$ is a probability vector. Let $d = (d_1, \dots, d_J)^T$ be a vector whose entries are nonnegative integers, with $K = \sum_{j=1}^J d_j$. Suppose that, for each j , p_j is the probability of index j and d_j is the number of times index j was chosen in K trials. The likelihood function of the parameters λ_i is

$$L(\lambda) = \prod_{j=1}^J p_j^{d_j} \quad (2.7)$$

so that the log-likelihood function is

$$LL(\lambda) = \sum_{j=1}^J d_j \log p_j. \quad (2.8)$$

Since p is a probability vector, maximizing $L(\lambda)$ is equivalent to minimizing $KL(d, p)$ with respect to λ , which, according to the lemma above, can be solved using SMART. In fact, since all of the block-iterative versions of SMART have the same limit whenever they have the same starting vector, any of these methods can be used to solve this maximum likelihood problem. In the case of transmission tomography the λ_i must be non-positive, so if SMART is to be used, some modification is needed to obtain such a solution.

We turn next to the block-iterative versions of the SMART, which we shall denote BI-SMART. These methods were known prior to the discovery of RBI-EMML and played an important role in that discovery; the importance of rescaling for acceleration was apparently not appreciated, however. The SMART was discovered in 1972, independently, by Darroch and Ratcliff [2], working in statistics, and by Schmidlin [3] in medical imaging. Block-iterative versions of SMART are also treated in [2], but they also insist on subset balance; the inconsistent case was not considered.

3 Block-iterative SMART

We start by considering a formulation of BI-SMART that is general enough to include all of the variants we wish to discuss. As we shall see, this formulation is too general and will need to be restricted in certain ways to obtain convergence. Let the BI-SMART iterative step be defined as

$$x_j^{k+1} = x_j^k \exp \left(\beta_{nj} \sum_{i \in B_n} \alpha_{ni} P_{ij} \log \left(\frac{y_i}{(Px^k)_i} \right) \right), \quad (3.1)$$

for $j = 1, 2, \dots, J$, $n = k(\bmod N) + 1$ and β_{nj} and α_{ni} positive. As we shall see, our convergence proof will require that β_{nj} be separable, that is,

$$\beta_{nj} = \gamma_j \delta_n$$

for each j and n so that (3.1) becomes

$$x_j^{k+1} = x_j^k \exp \left(\gamma_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \log \left(\frac{y_i}{(Px^k)_i} \right) \right). \quad (3.2)$$

We also require

$$\gamma_j \delta_n \sigma_{nj} \leq 1, \quad (3.3)$$

for $\sigma_{nj} = \sum_{i \in B_n} \alpha_{ni} P_{ij}$. With these conditions satisfied we have the following result.

Theorem 3.1 *Let there be nonnegative solutions of $y = Px$. For any positive vector x^0 and any collection of blocks $\{B_n, n = 1, \dots, N\}$ the sequence $\{x^k\}$ given by (3.2) converges to the unique solution of $y = Px$ for which the weighted cross-entropy $\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^0)$ is minimized.*

The inequality in the following lemma is the basis for the convergence proof.

Lemma 3.1 *Let $\beta_{nj} = \gamma_j \delta_n$ and $y = Px$ for some nonnegative x . Then for $\{x^k\}$ as in (3.2) we have*

$$\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^k) - \sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^{k+1}) \geq \delta_n \sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i). \quad (3.4)$$

Proof: Note that the quantity

$$\exp\left(\gamma_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \log\left(\frac{y_i}{(Px^k)_i}\right)\right)$$

in equation (3.2) can be written as

$$\exp\left((1 - \gamma_j \delta_n \sigma_{nj}) \log 1 + \gamma_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \log\left(\frac{y_i}{(Px^k)_i}\right)\right),$$

which, by the convexity of the exponential function, is not greater than

$$(1 - \gamma_j \delta_n \sigma_{nj}) + \gamma_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \frac{y_i}{(Px^k)_i}.$$

It follows that

$$\sum_{j=1}^J \gamma_j^{-1} (x_j^k - x_j^{k+1}) \geq \delta_n \sum_{i \in B_n} \alpha_{ni} ((Px^k)_i - y_i). \quad (3.5)$$

Note that it is at this step that we have used the separability of β_{nj} . We also have

$$\log(x_j^{k+1}/x_j^k) = \gamma_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \log \frac{y_i}{(Px^k)_i}. \quad (3.6)$$

Therefore

$$\begin{aligned} \sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^k) - \sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^{k+1}) &= \sum_{j=1}^J \gamma_j^{-1} (x_j \log(x_j^{k+1}/x_j^k) + x_j^k - x_j^{k+1}) \\ &= \sum_{j=1}^J x_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \log \frac{y_i}{(Px^k)_i} + \sum_{j=1}^J \gamma_j^{-1} (x_j^k - x_j^{k+1}) \\ &= \delta_n \sum_{i \in B_n} \alpha_{ni} \left(\sum_{j=1}^J x_j P_{ij} \right) \log \frac{y_i}{(Px^k)_i} + \sum_{j=1}^J \gamma_j^{-1} (x_j^k - x_j^{k+1}) \end{aligned}$$

$$\geq \delta_n \left(\sum_{i \in B_n} \alpha_{ni} (y_i \log \frac{y_i}{(Px^k)_i} + (Px^k)_i - y_i) \right) = \delta_n \sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i).$$

This completes the proof of the lemma. ■

From (3.4) the sequence $\{\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^k)\}$ is decreasing, from which we conclude both that its difference sequence, the left side of inequality (3.4), converges to zero and that the sequence $\{x^k\}$ is bounded. Since the left side of inequality (3.4) dominates the right side, the nonnegative sequence $\{\sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i)\}$ is also converging to zero. Let x^* be any cluster point of the sequence $\{x^k\}$. Then it is not difficult to show that $y = Px^*$. Replacing x with x^* we have that the sequence $\{\sum_{j=1}^J \gamma_j^{-1} KL(x_j^*, x_j^k)\}$ is decreasing; since a subsequence converges to zero, so does the whole sequence. Therefore x^* is the limit of the sequence $\{x^k\}$. This proves that the algorithm produces a solution of $y = Px$. To conclude further that the solution is the one for which the quantity $\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^0)$ is minimized requires further work to replace (3.4) with equation (5.10), in which the right side is independent of the particular solution x chosen; see the final section for the details.

We see from the theorem that how we select the γ_j is determined by how we wish to weight the terms in the sum $\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^0)$. In some cases we want to minimize the cross-entropy $KL(x, x^0)$ subject to $y = Px$; in this case we would select $\gamma_j = 1$. In other cases we may have some prior knowledge as to the relative sizes of the x_j and wish to emphasize the smaller values more; then we may choose γ_j proportional to our prior estimate of the size of x_j . Having selected the γ_j , we see from (3.4) that convergence will be accelerated if we select δ_n as large as permitted by the condition $\gamma_j \delta_n \sigma_{nj} \leq 1$. This suggests that we take

$$\delta_n = 1 / \max\{\sigma_{nj} \gamma_j, j = 1, \dots, J\}. \quad (3.7)$$

The *rescaled* BI-SMART (RBI-SMART) as presented in [13, 14, 15] uses this choice, but with $\alpha_{ni} = 1$ for each n and i . Let's look now at some of the other choices for these parameters that have been considered in the literature.

First, we notice that the OSSMART does not generally satisfy the requirements, since in (2.4) the choices are $\alpha_{ni} = 1$ and $\beta_{nj} = s_{nj}^{-1}$; the only times this is acceptable is if the s_{nj} are separable; that is, $s_{nj} = r_j t_n$ for some r_j and t_n . This is slightly more general than the condition of subset balance and is sufficient for convergence of OSSMART, since, for $\gamma_j = \alpha_{ni} = 1$ and δ_n as in (3.7), the BI-SMART reduces to the OSSMART.

In [4] Censor and Segman make the choices $\beta_{nj} = 1$ and $\alpha_{ni} > 0$ such that $\sigma_{nj} \leq 1$ for all n and j . In those cases in which σ_{nj} is much less than 1 for each n and j

their iterative scheme is probably excessively relaxed; it is hard to see how one might improve the rate of convergence by altering only the weights α_{ni} , however. Limiting the choice to $\gamma_j \delta_n = 1$ reduces our ability to accelerate this algorithm.

The original SMART in (2.1) uses $N = 1$, $\gamma_j = s_j^{-1}$ and $\alpha_{ni} = \alpha_i = 1$. Clearly (3.3) is satisfied; in fact it becomes an equality now.

For the row-action version of SMART, the *multiplicative* ART (MART), due to Gordon, Bender and Herman [11], we take $N = I$ and $B_n = B_i = \{i\}$ for $i = 1, \dots, I$. The MART begins with a strictly positive vector x^0 and has the iterative step

The MART:

$$x_j^{k+1} = x_j^k \left(\frac{y_i}{(Px^k)_i} \right)^{m_i^{-1} P_{ij}}, \quad (3.8)$$

for $j = 1, 2, \dots, J$, $i = k(\text{mod } I) + 1$ and $m_i > 0$ chosen so that $m_i^{-1} P_{ij} \leq 1$ for all j . Convergence of the MART is generally faster for smaller m_i , so a good choice is $m_i = \max\{P_{ij} | j = 1, \dots, J\}$. Although this particular choice for m_i is not explicitly mentioned in the various discussions of MART, it was used in implementations of MART from the beginning [16].

Darroch and Ratcliff included a discussion of a block-iterative version of SMART in their 1972 paper [2]. Close inspection of their version reveals that they require that $s_{nj} = \sum_{i \in B_n} P_{ij} = 1$ for all j . Since this is unlikely to be the case initially, we might try to rescale the equations or unknowns to obtain this condition. However, unless $s_{nj} = \sum_{i \in B_n} P_{ij}$ depends only on j and not on n , which is the *subset balance* property used in [12], we cannot redefine the unknowns in a way that is independent of n .

The MART fails to converge in the inconsistent case. What is always observed, but for which no proof exists, is that, for each fixed $i = 1, 2, \dots, I$, as $m \rightarrow +\infty$, the MART subsequences $\{x^{mI+i}\}$ converge to separate limit vectors, say $x^{\infty,i}$. This *limit cycle* $LC = \{x^{\infty,i} | i = 1, \dots, I\}$ reduces to a single vector whenever there is a nonnegative solution of $y = Px$. The greater the minimum value of $KL(Px, y)$ the more distinct from one another the vectors of the limit cycle are. An analogous result is observed for BI-SMART.

4 Block-iterative EMMML

As we did with SMART, we consider now a formulation of BI-EMML that is general enough to include all of the variants we wish to discuss. Once again, the formulation is too general and will need to be restricted in certain ways to obtain convergence.

Let the iterative step of the BI-EMML be defined as

$$x_j^{k+1} = x_j^k(1 - \beta_{nj}\sigma_{nj}) + x_j^k\beta_{nj} \sum_{i \in B_n} \alpha_{ni}P_{ij} \frac{y_i}{(Px^k)_i}, \quad (4.1)$$

for $j = 1, 2, \dots, J$, $n = k(\bmod N) + 1$ and β_{nj} and α_{ni} positive. As in the case of BI-SMART, our convergence proof will require that β_{nj} be separable, that is,

$$\beta_{nj} = \gamma_j\delta_n$$

for each j and n and that (3.3) hold. The BI-EMML then becomes

$$x_j^{k+1} = x_j^k(1 - \gamma_j\delta_n\sigma_{nj}) + x_j^k\gamma_j\delta_n \sum_{i \in B_n} \alpha_{ni}P_{ij} \frac{y_i}{(Px^k)_i}, \quad (4.2)$$

With these conditions satisfied we have the following result.

Theorem 4.1 *Let there be nonnegative solutions of $y = Px$. For any positive vector x^0 and any collection of blocks $\{B_n, n = 1, \dots, N\}$ the sequence $\{x^k\}$ given by (4.2) converges to a nonnegative solution of $y = Px$.*

When there are multiple nonnegative solutions of $y = Px$ the solution obtained by BI-EMML will depend on the starting point x^0 , but precisely how it depends on x^0 is an open question. Also, in contrast to the case of BI-SMART, the solution can depend on the particular choice of the blocks. The inequality in the following lemma is the basis for the convergence proof.

Lemma 4.1 *Let $y = Px$ for some nonnegative x . Then for $\{x^k\}$ as in (4.2) we have*

$$\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^k) - \sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^{k+1}) \geq \delta_n \sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i). \quad (4.3)$$

Proof: From the iterative step

$$x_j^{k+1} = x_j^k(1 - \gamma_j\delta_n\sigma_{nj}) + x_j^k\gamma_j\delta_n \sum_{i \in B_n} \alpha_{ni}P_{ij} \frac{y_i}{(Px^k)_i}$$

we have

$$\log(x_j^{k+1}/x_j^k) = \log\left((1 - \gamma_j\delta_n\sigma_{nj}) + \gamma_j\delta_n \sum_{i \in B_n} \alpha_{ni}P_{ij} \frac{y_i}{(Px^k)_i}\right).$$

By the concavity of the logarithm we obtain the inequality

$$\log(x_j^{k+1}/x_j^k) \geq \left((1 - \gamma_j\delta_n\sigma_{nj}) \log 1 + \gamma_j\delta_n \sum_{i \in B_n} \alpha_{ni}P_{ij} \log \frac{y_i}{(Px^k)_i}\right),$$

or

$$\log(x_j^{k+1}/x_j^k) \geq \gamma_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \log \frac{y_i}{(Px^k)_i}.$$

Therefore

$$\sum_{j=1}^J \gamma_j^{-1} x_j \log(x_j^{k+1}/x_j^k) \geq \delta_n \sum_{i \in B_n} \alpha_{ni} \left(\sum_{j=1}^J x_j P_{ij} \right) \log \frac{y_i}{(Px^k)_i}. \quad (4.4)$$

Note that it is at this step that we used the separability of the β_{nj} . Also

$$\sum_{j=1}^J \gamma_j^{-1} (x_j^{k+1} - x_j^k) = \delta_n \sum_{i \in B_n} ((Px^k)_i - y_i). \quad (4.5)$$

Since the left sides and right sides of inequalities (4.4) and (4.5) add to the left side and right side of inequality (4.3), respectively, this concludes the proof of the lemma. \blacksquare

From (4.3) we conclude, as we did in the BI-SMART case, that the sequence $\{\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^k)\}$ is decreasing, that $\{x^k\}$ is therefore bounded and the sequence $\{\sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i)\}$ is converging to zero. Let x^* be any cluster point of the sequence $\{x^k\}$. Then it is not difficult to show that $y = Px^*$. Replacing x with x^* we have that the sequence $\{\sum_{j=1}^J \gamma_j^{-1} KL(x_j^*, x_j^k)\}$ is decreasing; since a subsequence converges to zero, so does the whole sequence. Therefore x^* is the limit of the sequence $\{x^k\}$. This proves that the algorithm produces a nonnegative solution of $y = Px$. We have been unable to replace (4.3) with an equation in which the right side is independent of the particular solution x chosen; for that reason we can say no more about the solution that has been obtained.

Having selected the γ_j , we see from (4.3) that convergence will be accelerated if we select δ_n as large as permitted by the condition $\gamma_j \delta_n \sigma_{nj} \leq 1$. This suggests that once again we take δ_n as in (3.7). The *rescaled* BI-EMML (RBI-EMML) as presented in [13, 14, 15] uses this choice, but with $\alpha_{ni} = 1$ for each n and i . Let's look now at some of the other choices for these parameters that have been considered in the literature.

First, we notice that the OSEM does not generally satisfy the requirements, since in (2.3) the choices are $\alpha_{ni} = 1$ and $\beta_{nj} = s_{nj}^{-1}$; the only times this is acceptable is if the s_{nj} are separable; that is, $s_{nj} = r_j t_n$ for some r_j and t_n . This is slightly more general than the condition of subset balance and is sufficient for convergence of OSEM, since, for $\gamma_j = \alpha_{ni} = 1$ and δ_n as in (3.7), the BI-EMML reduces to the OSEM.

The original EMML in (2.2) uses $N = 1$, $\gamma_j = s_j^{-1}$ and $\alpha_{ni} = \alpha_i = 1$. Clearly (3.3) is satisfied; in fact it becomes an equality now.

Notice that the calculations required to perform the BI-SMART are somewhat more complicated than those needed in BI-EMML. Because the MART converges rapidly in most cases there is considerable interest in the row-action version of EMML. It was clear from the outset that using the OSEM in a row-action mode does not work. We see from the formula for BI-EMML that the proper row-action version of EMML, which we call the EM-MART, has the iterative step

EM-MART:

$$x_j^{k+1} = (1 - \delta_i \gamma_j \alpha_{ii} P_{ij}) x_j^k + \delta_i \gamma_j \alpha_{ii} P_{ij} \frac{y_i}{(Px^k)_i}, \quad (4.6)$$

with

$$\gamma_j \delta_i \alpha_{ii} P_{ij} \leq 1$$

for all i and j . The optimal choice would seem to be to take $\delta_i \alpha_{ii}$ as large as possible; that is, to select $\delta_i \alpha_{ii} = 1 / \max\{\gamma_j P_{ij}, j = 1, \dots, J\}$. With this choice the EM-MART is called the *rescaled* EM-MART (REM-MART).

The EM-MART fails to converge in the inconsistent case. What is always observed, but for which no proof exists, is that, for each fixed $i = 1, 2, \dots, I$, as $m \rightarrow +\infty$, the EM-MART subsequences $\{x^{mI+i}\}$ converge to separate limit vectors, say $x^{\infty,i}$. This *limit cycle* $LC = \{x^{\infty,i} | i = 1, \dots, I\}$ reduces to a single vector whenever there is a nonnegative solution of $y = Px$. The greater the minimum value of $KL(y, Px)$ the more distinct from one another the vectors of the limit cycle are. An analogous result is observed for BI-EMML.

We must mention a method that closely resembles the REM-MART, the *row-action maximum likelihood algorithm* (RAMLA), which was discovered independently by Browne and De Pierro [17]. The RAMLA avoids the limit cycle in the inconsistent case by using strong underrelaxation involving a decreasing sequence of relaxation parameters λ_k . The RAMLA has the following iterative step:

RAMLA:

$$x_j^{k+1} = (1 - \lambda_k \sum_{i \in B_n} P_{ij}) x_j^k + \lambda_k x_j^k \sum_{i \in B_n} P_{ij} \left(\frac{y_i}{(Px^k)_i} \right), \quad (4.7)$$

where the positive relaxation parameters λ_k are chosen to converge to zero and $\sum_{k=0}^{+\infty} \lambda_k = +\infty$.

5 Proof of convergence of BI-SMART

As we stated earlier, in the consistent case the sequence $\{x^k\}$ generated by the BI-SMART algorithm and given by equation (3.2) converges to the unique solution of

$y = Px$ for which the distance $\sum_{j=1}^J \gamma_j^{-1} KL(x_j, x_j^0)$ is minimized. In this section we sketch the proof of this result through a sequence of lemmas, each of which is easily established.

Lemma 5.1 *For any nonnegative vectors a and b with $a_+ = \sum_{m=1}^M a_m$ and $b_+ = \sum_{m=1}^M b_m > 0$ we have*

$$KL(a, b) = KL(a_+, b_+) + KL(a, \frac{a_+}{b_+} b). \quad (5.1)$$

so that $KL(a, b) \geq KL(a_+, b_+)$.

For nonnegative vectors x and z let

$$G_n(x, z) = \sum_{j=1}^J \gamma_j^{-1} KL(x_j, z_j) + \delta_n \sum_{i \in B_n} \alpha_{ni} [KL((Px)_i, y_i) - KL((Px)_i, (Pz)_i)]. \quad (5.2)$$

It follows from (5.1) and the inequality

$$\gamma_j^{-1} - \delta_n \sigma_{nj} \geq 1$$

that

$$\sum_{j=1}^J \gamma_j^{-1} KL(x_j, z_j) - \delta_n \sum_{i \in B_n} \alpha_{ni} KL((Px)_i, (Pz)_i) \geq 0 \quad (5.3)$$

and so $G_n(x, z) \geq 0$ in all cases.

Lemma 5.2 *For every x we have*

$$G_n(x, x) = \delta_n \sum_{i \in B_n} \alpha_{ni} KL((Px)_i, y_i) \quad (5.4)$$

so that

$$G_n(x, z) = G_n(x, x) + \sum_{j=1}^J \gamma_j^{-1} KL(x_j, z_j) - \delta_n \sum_{i \in B_n} \alpha_{ni} KL((Px)_i, (Pz)_i). \quad (5.5)$$

Therefore the distance $G_n(x, z)$ is minimized, as a function of z , by $z = x$.

Now we minimize $G_n(x, z)$ as a function of x .

Lemma 5.3 *For each x and z we have*

$$G_n(x, z) = G_n(z', z) + \sum_{j=1}^J \gamma_j^{-1} KL(x_j, z'_j), \quad (5.6)$$

where

$$z'_j = z_j \exp \left(\gamma_j \delta_n \sum_{i \in B_n} \alpha_{ni} P_{ij} \log \frac{y_i}{(Pz)_i} \right) \quad (5.7)$$

for each z .

It is clear that $(x^k)' = x^{k+1}$ for all k ; this lemma motivates the definition of the iterative step in the BI-SMART.

Now let $y = Pu$ for some nonnegative vector u . We calculate $G_n(u, x^k)$ in two ways: using the definition we have

$$G_n(u, x^k) = \sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^k) - \delta_n \sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i), \quad (5.8)$$

while using (5.6) we find that

$$G_n(u, x^k) = G_n(x^{k+1}, x^k) + \sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^{k+1}). \quad (5.9)$$

Therefore

$$\sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^k) - \sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^{k+1}) = G_n(x^{k+1}, x^k) + \delta_n \sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i). \quad (5.10)$$

We conclude several things from this.

First, the sequence $\{\sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^k)\}$ is decreasing, so that the sequences $\{G_n(x^{k+1}, x^k)\}$ and $\{\delta_n \sum_{i \in B_n} \alpha_{ni} KL(y_i, (Px^k)_i)\}$ converge to zero. Therefore the sequence $\{x^k\}$ is bounded and we may select an arbitrary cluster point x^* . It follows that $y = Px^*$. We may therefore replace the generic solution u with x^* to find that the sequence $\{\sum_{j=1}^J \gamma_j^{-1} KL(x_j^*, x_j^k)\}$ is decreasing; but since a subsequence is converging to zero, the entire sequence must converge to zero. Therefore $\{x^k\}$ converges to the solution x^* .

Finally, since the right side of (5.10) does not depend on the particular choice of solution we have made, neither does the left side. By *telescoping*, that is, by summing on k on both sides, we conclude that

$$\sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^0) - \sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^*)$$

is also independent of the choice of u . Consequently, minimizing $\sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^0)$ over all solutions u is equivalent to minimizing $\sum_{j=1}^J \gamma_j^{-1} KL(u_j, x_j^*)$ over all solutions u ; but the solution to the latter problem is obviously $u = x^*$. This completes the proof.

Acknowledgments

This work was supported in part by grants R01EB001457 and R01EB002798 from the National Institute of Biomedical Imaging and Bioengineering. The views expressed here are those of the author and do not necessarily reflect those of the funder.

References

- [1] Y. Vardi, L.A. Shepp and L. Kaufman, A statistical model for positron emission tomography, *Journal of the American Statistical Association*, **80**, pp. 8–20, 1985.
- [2] J. Darroch and D. Ratcliff, Generalized iterative scaling for log-linear models, *The Annals of Mathematical Statistics*, **43** (5), pp. 1470–1480, 1972.
- [3] P. Schmidlin, Iterative separation of sections in tomographic scintigrams, *Nuclear Medicine*, **15** (1), Schatten Verlag, Stuttgart, 1972.
- [4] Y. Censor and J. Segman, On block-iterative maximization, *Journal of Information and Optimization Sciences*, **8**, pp. 275–291, 1987.
- [5] C. Byrne, Iterative image reconstruction algorithms based on cross-entropy minimization, *IEEE Transactions on Image Processing*, **IP-2**, pp. 96–103, 1993.
- [6] C. Byrne, Erratum and addendum to “Iterative image reconstruction algorithms based on cross-entropy minimization”, *IEEE Transactions on Image Processing*, **IP-4**, pp. 225–226, 1995.
- [7] R. Leahy and C. Byrne, Guest editorial: Recent development in iterative image reconstruction for PET and SPECT, *IEEE Transactions on Medical Imaging*, **19**, pp. 257–260, 2000.
- [8] S. Kullback and R. Leibler, On information and sufficiency, *Annals of Mathematical Statistics*, **22**, pp. 79–86, 1951.
- [9] I. Csiszár, A geometric interpretation of Darroch and Ratcliff’s generalized iterative scaling, *The Annals of Statistics*, **17** (3), pp. 1409–1413, 1989.
- [10] I. Csiszár, Why least squares and maximum entropy? An axiomatic approach to inference for linear inverse problems, *The Annals of Statistics*, **19** (4), pp. 2032–2066, 1991.

- [11] R. Gordon, R. Bender and G.T. Herman, Algebraic reconstruction techniques (ART) for three-dimensional electron microscopy and x-ray photography, *Journal of Theoretical Biology*, **29**, pp. 471–481, 1970.
- [12] H.M. Hudson and R.S. Larkin, Accelerated image reconstruction using ordered subsets of projection data, *IEEE Transactions on Medical Imaging*, **13**, pp. 601–609, 1994.
- [13] C. Byrne, Block-iterative methods for image reconstruction from projections, *IEEE Transactions on Image Processing*, **IP-5**, pp. 792–794, 1996.
- [14] C. Byrne, Convergent block-iterative algorithms for image reconstruction from inconsistent data, *IEEE Transactions on Image Processing*, **IP-6**, pp. 1296–1304, 1997.
- [15] C. Byrne, Accelerating the EMLL algorithm and related iterative algorithms by rescaled block-iterative methods, *IEEE Transactions on Image Processing*, **IP-7**, pp. 100–109, 1998.
- [16] G.T. Herman, *private communication*, 1999.
- [17] J. Browne and A. De Pierro, A row-action alternative to the EM algorithm for maximizing likelihoods in emission tomography, *IEEE Transactions on Medical Imaging*, **15**, pp. 687–699, 1996.