

On complete-data spaces for PET reconstruction algorithms

Jeffrey A. Fessler, Neal H. Clinthorne, and W. Leslie Rogers
University of Michigan

ABSTRACT

As investigators consider more comprehensive measurement models for emission tomography, there will be more choices for the complete-data spaces of the associated expectation-maximization (EM) algorithms for maximum-likelihood (ML) estimation. In this paper, we show that EM algorithms based on smaller complete-data spaces will typically converge faster. We discuss two practical applications of these concepts: (i) the ML-IA and ML-IB image reconstruction algorithms of Politte and Snyder [1] which are based on measurement models that account for attenuation and accidental coincidences in positron-emission tomography (PET), and (ii) the problem of simultaneous estimation of emission and transmission parameters. Although the PET applications may often violate the necessary regularity conditions, our analysis predicts heuristically that the ML-IB algorithm, which has a smaller complete-data space, should converge faster than ML-IA. This is corroborated by the empirical findings in [1].

I. INTRODUCTION

The ML criterion for tomographic image reconstruction has received considerable attention since Shepp and Vardi [2] introduced an EM algorithm for computing ML estimates. Although the medical imaging community often refers to "the" ML-EM algorithm, there are in fact a multitude of feasible EM algorithms, each based on a different complete-data space. A useful complete-data space supplements the observed measurements in a way that facilitates parameter estimation [3]. Although only one complete-data space has been suggested for PET under the simple measurement model used in the early papers [2, 4], there will be more choices as investigators consider more comprehensive measurement models, such as those accounting for photon attenuation [5], accidental coincidences [1], deadtime, and scatter [6, 7]. This paper illustrates the importance of parsimony in choosing complete-data spaces, and some of the tradeoffs that result.

Accurate quantification of radiotracer activity using PET must include corrections for the effects of attenuation and accidental coincidences. Recently, Politte and Snyder proposed two ML-EM algorithms for PET image reconstruction that directly incorporate the effects of known attenuation and accidental coincidences into the statistical

This work was supported in part by a National Institute of Health Training Fellowship, a DOE Alexander Hollaender Postdoctoral Fellowship, and DOE Grant DE-FG02-87ER65061.

measurement model [1]. The algorithms are based on two different complete-data spaces, one of which is a subset of the other. They observed in experiments that the algorithm based on the smaller complete-data space converged faster. In this paper we corroborate their observations by proving that smaller complete-data spaces yield EM algorithms with faster asymptotic convergence rates.

The measurement models used in [1] assumed exact knowledge of the survival probabilities, the probability that both photons of a positron-produced pair escape unattenuated. In practice, one must obtain these factors experimentally, typically by a transmission scan that precedes the radiotracer injection. As mentioned in [1], a more accurate approach would account for the statistical uncertainties in both the emission data and the transmission measurements. An iterative method for simultaneously estimating the emission intensities and the survival probabilities has been recently proposed by Clinthorne *et al.* [8]. In this paper, we present two algorithms for joint emission/transmission estimation based on generalizations of the two complete-data spaces in [1]. We demonstrate that although the smaller complete-data space may provide a faster algorithm in theory, in practice the larger complete-data space leads to an EM algorithm with an easier maximization step. Such tradeoffs have been observed in other EM applications [3], but they may be particularly important to future investigations of PET reconstruction methods because of the large dimensions of the parameter spaces.

In Section II, we briefly review the EM algorithm, and prove that smaller complete-data spaces result in faster convergence. This is applied to PET in Section III. In Section IV we analyze the joint emission/transmission estimation method for PET.

II. THEORY

A. EM algorithm

We observe \mathbf{y} , a realization of a random vector \mathbf{Y} having known density $g(\mathbf{y}; \boldsymbol{\theta})$, with the goal of computing the ML estimate of $\boldsymbol{\theta}$. In many problems, including emission tomography, the measurements are "incomplete" in the sense that several components of $\boldsymbol{\theta}$ may contribute to each component of \mathbf{Y} . In such problems one can often postulate a "complete data" random vector \mathbf{X} that is more naturally related to the parameter vector $\boldsymbol{\theta}$, and is related to the observed measurements by a many-to-one mapping $\mathbf{Y} = h(\mathbf{X})$. The density $f(\mathbf{x}; \boldsymbol{\theta})$ of the complete data \mathbf{X} must

be consistent with the incomplete data \mathbf{Y} in that:

$$g(\mathbf{y}; \theta) = \int_{\{\mathbf{x}:\mathbf{y}=h(\mathbf{x})\}} f(\mathbf{x}; \theta) d\mathbf{x}.$$

Let

$$Q(\bar{\theta}; \theta) \triangleq E \{ \log f(\mathbf{X}; \bar{\theta}) | \mathbf{Y} = \mathbf{y}; \theta \} \quad (1)$$

$$\begin{aligned} &= \int \log f(\mathbf{x}; \bar{\theta}) f(\mathbf{x} | \mathbf{Y} = \mathbf{y}; \theta) d\mathbf{x} \\ &= H(\bar{\theta}; \theta) + L(\bar{\theta}), \end{aligned} \quad (2)$$

where

$$\begin{aligned} H(\bar{\theta}; \theta) &\triangleq E \{ \log f(\mathbf{X} | \mathbf{Y} = \mathbf{y}; \bar{\theta}) | \mathbf{Y} = \mathbf{y}; \theta \}, \\ L(\theta) &\triangleq \log g(\mathbf{y}; \theta). \end{aligned}$$

The EM algorithm for ML estimation [3], calls for iterating over the following steps:

E-step:

$$\text{Compute } Q(\theta; \theta^i),$$

M-step:

$$\theta^{i+1} = \arg \max_{\theta} Q(\theta; \theta^i),$$

where θ^i denotes the parameter estimate after the i th iteration. Note that by Jensen's inequality [3]:

$$H(\bar{\theta}; \theta) \leq H(\theta; \theta) \quad \forall \bar{\theta},$$

so the EM algorithm produces a likelihood sequence $L(\theta^i)$ that is monotonically increasing. The basic idea is to compute Q , the conditional expectation of the complete data given the most recent parameter estimate, and then to maximize the parameter's likelihood as if one had observed the complete data [3]. The EM algorithm is most useful when the complete-data space is chosen such that $Q(\theta; \theta^i)$ can be maximized analytically for the M-step, although other approaches are possible [9].

B. EM convergence rate

Several investigators have observed empirically that larger complete-data spaces correspond to slower EM convergence [3, pp. 25,34]. In this section we formally establish a version of this result. For our purposes, asymptotic convergence rate is defined by the following well known result [10, p. 301].

Linear Convergence Theorem: *If (i) $G : \mathcal{D} \in \mathbb{R}^n \rightarrow \mathbb{R}^n$ has a fixed point $\theta^* \in \mathcal{D}_+ = \text{int}(\mathcal{D})$, (ii) G is Fréchet differentiable at θ^* , and (iii) $\rho(\nabla^1 G(\theta^*)) < 1$, where $\rho(\cdot)$ denotes spectral radius¹, then the root-convergence factor [10, p.*

¹The (n, m) element of the $\text{dim}(\theta)$ by $\text{dim}(\theta)$ matrix $\nabla^1 G$ is $\partial/\partial\theta_m G_n(\theta)$. In general, the $\text{dim}(\theta)$ by $\text{dim}(\theta)$ matrix $\nabla^{ij} F(\bar{\theta}; \theta)$ denotes the matrix of partials $\partial^{i+j} F(\bar{\theta}; \theta) / \partial \bar{\theta}^i \partial \theta^j$.

288] R_1 at θ^* for the iterative process $\theta^{i+1} = G(\theta^i)$ is given by $R_1(G, \theta^*) = \rho(\nabla^1 G(\theta^*))$.

This leads to the following definition of convergence rate for EM algorithms corresponding to strictly concave likelihood functions [3].

Theorem 1: *Let $\theta^{i+1} = G(\theta^i)$ define the iterations for an EM algorithm such that (i) G and θ^* satisfy conditions (i) and (ii) of the Linear Convergence Theorem, (ii) G is defined by solving the system of equations $\nabla^{10} Q(\bar{\theta}; \theta) |_{\bar{\theta}=G(\theta)} = \mathbf{0}$, and (iii) $\mathbf{L} \triangleq -\nabla^2 L(\theta^*)$ is positive definite, then the root-convergence factor at θ^* for the EM iteration G is*

$$R_1 = \rho(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{L}) < 1, \quad (3)$$

where \mathbf{I} is the $n \times n$ identity matrix and $\mathbf{Q} \triangleq -\nabla^{20} Q(\theta^*; \theta^*)$.

Proof: $\mathbf{H} \triangleq -\nabla^{20} H(\theta^*; \theta^*)$ is nonnegative definite since it is a (conditional) Fisher information matrix [11, p. 126]. From (2), $\mathbf{Q} = \mathbf{H} + \mathbf{L}$, so by (iii), \mathbf{Q} is positive definite and therefore invertible. From (ii) we see that $\nabla^{10} Q(G(\theta); \theta) = \mathbf{0}$, which can be differentiated again to yield $\nabla^{20} Q(G(\theta); \theta) \nabla^1 G(\theta) + \nabla^{11} Q(G(\theta); \theta) = \mathbf{0}$. Therefore,

$$\nabla^1 G(\theta^*) = -(\nabla^{20} Q(\theta^*; \theta^*))^{-1} \nabla^{11} Q(\theta^*; \theta^*),$$

since $G(\theta^*) = \theta^*$ by (i). But it is well known [3] that for $\theta \in \mathcal{D}_+$

$$\begin{aligned} \nabla^{11} Q(\theta; \theta) &= \nabla^{11} H(\theta; \theta) = -\nabla^{20} H(\theta; \theta) \\ &= \nabla^2 L(\theta) - \nabla^{20} Q(\theta; \theta). \end{aligned}$$

The theorem then follows from combining the two equations above and the Linear Convergence Theorem, provided $\rho(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{L}) < 1$. To prove this, we parallel an argument of Green [12]. If α is an eigenvalue of $\mathbf{I} - \mathbf{Q}^{-1}\mathbf{L}$, then $|\mathbf{I} - \mathbf{Q}^{-1}\mathbf{L} - \alpha\mathbf{I}| = 0$, hence $|(1 - \alpha)\mathbf{Q} - \mathbf{L}| = 0$. Since $\mathbf{Q} = \mathbf{H} + \mathbf{L}$, $|(1 - \alpha)\mathbf{H} - \alpha\mathbf{L}| = 0$. Thus, by assumption (iii), we must have $\alpha \in [0, 1)$, hence $\rho(\mathbf{I} - \mathbf{Q}^{-1}\mathbf{L}) < 1$. \square

Since $\mathbf{Q} = \mathbf{H} + \mathbf{L}$, where \mathbf{H} is a conditional Fisher information matrix, one sees from (3) that if a larger complete-data space has more Fisher information, then the corresponding root-convergence factor will be larger, and the asymptotic convergence rate will be slower. This is the idea behind the next lemma and theorem, the main results of this section.

Lemma 1: *If (i) $\mathbf{Q}_B = \mathbf{H} + \mathbf{L}$, where \mathbf{H} is symmetric nonnegative definite and \mathbf{L} is symmetric positive definite, and (ii) $\mathbf{Q}_A = \mathbf{Q}_B + \mathbf{N}$ where \mathbf{N} is symmetric nonnegative definite, then $\rho_B \leq \rho_A$, where $\rho_A = \rho(\mathbf{I} - \mathbf{Q}_A^{-1}\mathbf{L})$ and $\rho_B = \rho(\mathbf{I} - \mathbf{Q}_B^{-1}\mathbf{L})$. Furthermore, if \mathbf{N} is symmetric positive definite, then $\rho_B < \rho_A$.*

Proof: Again we borrow from Green [12]. By the arguments in Theorem 1, $0 \leq \rho_A < 1$ and $0 \leq \rho_B < 1$. Since

$\rho_B = \rho(\mathbf{I} - \mathbf{Q}_B^{-1}\mathbf{L})$, $\exists \mathbf{u} \neq 0$ s.t. $(\mathbf{I} - \mathbf{Q}_B^{-1}\mathbf{L})\mathbf{u} = \rho_B\mathbf{u}$, so $(1 - \rho_B)\mathbf{Q}_B\mathbf{u} - \mathbf{L}\mathbf{u} = 0$. By (ii), $(1 - \rho_B)\mathbf{Q}_A\mathbf{u} - \mathbf{L}\mathbf{u} = (1 - \rho_B)\mathbf{N}\mathbf{u}$, so $(1 - \rho_B)\mathbf{Q}_A^{\frac{1}{2}}\mathbf{u} - \mathbf{Q}_A^{-\frac{1}{2}}\mathbf{L}\mathbf{u} = (1 - \rho_B)\mathbf{Q}_A^{-\frac{1}{2}}\mathbf{N}\mathbf{u}$. Defining $\mathbf{v} = \mathbf{Q}_A^{\frac{1}{2}}\mathbf{u}$, it follows that $(1 - \rho_B)\mathbf{v}'\mathbf{v} - \mathbf{v}'\mathbf{Q}_A^{-\frac{1}{2}}\mathbf{L}\mathbf{Q}_A^{-\frac{1}{2}}\mathbf{v} = (1 - \rho_B)\mathbf{v}'\mathbf{Q}_A^{-\frac{1}{2}}\mathbf{N}\mathbf{Q}_A^{-\frac{1}{2}}\mathbf{v} \geq 0$, if \mathbf{N} is nonnegative definite. Hence, $\mathbf{v}'[\mathbf{I} - \mathbf{Q}_A^{-\frac{1}{2}}\mathbf{L}\mathbf{Q}_A^{-\frac{1}{2}}]\mathbf{v} \geq \rho_B\mathbf{v}'\mathbf{v}$, so $\rho(\mathbf{I} - \mathbf{Q}_A^{-\frac{1}{2}}\mathbf{L}\mathbf{Q}_A^{-\frac{1}{2}}) \geq \rho_B$. But $\rho(\mathbf{I} - \mathbf{Q}_A^{-\frac{1}{2}}\mathbf{L}\mathbf{Q}_A^{-\frac{1}{2}}) = \rho(\mathbf{I} - \mathbf{Q}_A^{-1}\mathbf{L}) = \rho_A$, so $\rho_A \geq \rho_B$. The case when \mathbf{N} is positive definite is similar. \square

Theorem 2: *If (i) G_A and G_B are two EM algorithms that satisfy the conditions of Theorem 1 and that correspond to complete-data spaces \mathbf{X}_A and \mathbf{X}_B respectively, (ii) \mathbf{X}_B is a subset of \mathbf{X}_A , i.e., $\mathbf{X}_A = [\mathbf{X}_B, \mathbf{X}_o]'$, (iii) $f_A(\mathbf{x}_A|\mathbf{y}; \theta) = f_A([\mathbf{x}_B, \mathbf{x}_o]|\mathbf{y}; \theta) = f_B(\mathbf{x}_B|\mathbf{y}; \theta)f_o(\mathbf{x}_o|\mathbf{y}; \theta)$, and (iv) $f_o(\mathbf{x}_o|\mathbf{y}; \theta) = f_o(\mathbf{x}_o; \theta)$, i.e., \mathbf{X}_o is extraneous complete-data, then algorithm B converges faster than algorithm A asymptotically at a common fixed point θ^* .*

Proof: By integrating (iii) and using (iv), one sees that $f_A(\mathbf{x}_A; \theta) = f_A([\mathbf{x}_B, \mathbf{x}_o]; \theta) = f_B(\mathbf{x}_B; \theta)f_o(\mathbf{x}_o; \theta)$. Let Q_A and Q_B denote the Q function (1) for \mathbf{X}_A and \mathbf{X}_B respectively, then

$$\begin{aligned} Q_A(\bar{\theta}; \theta) &= \int \log f_A(\mathbf{x}_A; \bar{\theta}) f_A(\mathbf{x}_A|\mathbf{y}; \theta) d\mathbf{x}_A \\ &= \int \log f_A([\mathbf{x}_B, \mathbf{x}_o]; \bar{\theta}) f_A([\mathbf{x}_B, \mathbf{x}_o]|\mathbf{y}; \theta) d\mathbf{x}_B d\mathbf{x}_o \\ &= \int \log(f_B(\mathbf{x}_B; \bar{\theta})f_o(\mathbf{x}_o; \bar{\theta})) f_B(\mathbf{x}_B|\mathbf{y}; \theta)f_o(\mathbf{x}_o|\mathbf{y}; \theta) d\mathbf{x}_B d\mathbf{x}_o \\ &= \int \log f_B(\mathbf{x}_B; \bar{\theta}) f_B(\mathbf{x}_B|\mathbf{y}; \theta) d\mathbf{x}_B \\ &\quad + \int \log f_o(\mathbf{x}_o; \bar{\theta}) f_o(\mathbf{x}_o; \theta) d\mathbf{x}_o \\ &= Q_B(\bar{\theta}; \theta) + E\{\log f_o(\mathbf{X}_o; \bar{\theta}); \theta\}. \end{aligned}$$

Therefore,

$$-\nabla^{20}Q_A(\theta; \theta) = -\nabla^{20}Q_B(\theta; \theta) + \mathbf{J}_o(\theta),$$

where $\mathbf{J}_o(\theta)$ is the Fisher information matrix for the extraneous complete data \mathbf{X}_o at θ . The conclusion then follows from Lemma 1 and Theorem 1. \square

In summary, Theorem 2 shows that a complete-data space with extraneous variables will lead to an EM algorithm with slower convergence². Our derivation of this theorem relies on several assumptions, including convergence to an interior point, and strict concavity of the likelihood. The applicability of these assumptions in the PET context is discussed in the next section.

²We have since shown in [14] that Theorem 2 is true under considerably less restrictive conditions than (ii)-(iv), but this version is sufficient for the purposes of this paper.

III. PET RECONSTRUCTION

A. Two complete-data spaces

In this section, we briefly review the two PET reconstruction algorithms investigated in [1], and then discuss their convergence rates. For simplicity, we assume that the object is discretized into B voxels, and denote the radioactivity in the b th voxel by λ_b . PET measurements are acquired using D detector pairs, with the number of counts recorded by the d th pair denoted Y_d . The iterative algorithms attempt to estimate $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_B]'$ (corresponding to the parameter θ in the preceding section) from a realization \mathbf{y} of $\mathbf{Y} = [Y_1, \dots, Y_D]'$. The parameter domain is $\mathcal{D} = \{\boldsymbol{\lambda} : \lambda_b \geq 0, b = 1, \dots, B\}$.

Let p_{db} denote the point-spread function of the system, normalized so that $\sum_{d=1}^D p_{db} = 1$. Let q_b denote the detection probability for an unattenuated event originating in voxel b , and α_d denote the survival probability, assumed known here, for a photon pair emitted towards the d th detector pair. Then $\alpha_d w_{db}$ is the probability that an event in voxel b is detected by the d th detector pair, where $w_{db} = p_{db}q_b$. Let N_{db} denote the number of events from voxel b contributing to detector pair d ; the N_{db} 's have independent Poisson distributions with means $\alpha_d w_{db} \lambda_b$. Let R_d denote the number of accidental coincidences counted by detector pair d ; the R_d 's have independent Poisson distributions with mean r_d , assumed known. The total counts in the d th detector pair is then

$$Y_d = \sum_{b=1}^B N_{db} + R_d, \quad d = 1, \dots, D. \quad (4)$$

The two EM algorithms described in [1] were called ML-IA and ML-IB. The complete-data space for ML-IB is

$$\mathbf{X}_{IB} = \{\{N_{db}\}, \{R_d\}\}, \quad d = 1, \dots, D, \quad b = 1, \dots, B.$$

Under the distributions given above,

$$\begin{aligned} \log f_{IB}(\mathbf{X}_{IB}; \bar{\boldsymbol{\lambda}}) &= \sum_{d=1}^D (-r_d + R_d \log(r_d)) \\ &\quad + \sum_{d=1}^D \sum_{b=1}^B (-\alpha_d w_{db} \bar{\lambda}_b + N_{db} \log(\alpha_d w_{db} \bar{\lambda}_b)). \end{aligned} \quad (5)$$

It follows from (4) that [1]

$$E\{N_{db}|\mathbf{y}; \boldsymbol{\lambda}\} = y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\boldsymbol{\lambda})} \quad (6)$$

$$E\{R_d|\mathbf{y}; \boldsymbol{\lambda}\} = y_d \frac{r_d}{\hat{y}_d(\boldsymbol{\lambda})}, \quad (7)$$

where

$$\hat{y}_d(\boldsymbol{\lambda}) = \sum_{b=1}^B \alpha_d w_{db} \lambda_b + r_d. \quad (8)$$

Combining equations (5)-(7) and (1) yields

$$Q_{IB}(\bar{\lambda}; \lambda) = \sum_{d=1}^D \left(-r_d + y_d \frac{r_d}{\hat{y}_d(\lambda)} \log(r_d) \right) + \sum_{d=1}^D \sum_{b=1}^B \left(-\alpha_d w_{db} \bar{\lambda}_b + y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\lambda)} \log(\alpha_d w_{db} \bar{\lambda}_b) \right). \quad (9)$$

Maximizing $Q_{IB}(\bar{\lambda}; \lambda)$ over $\bar{\lambda}$ yields the ML-IB iteration $\lambda^{i+1} = G_{IB}(\lambda^i)$, where

$$G_{IB}(\lambda) \triangleq \lambda \odot [\mathbf{W}'(\alpha \odot \mathbf{y} \odot \hat{\mathbf{y}}(\lambda))] \odot [\mathbf{W}'\alpha],$$

$\alpha = [\alpha_1, \dots, \alpha_D]'$, $\hat{\mathbf{y}}(\lambda) = [\hat{y}_1(\lambda), \dots, \hat{y}_D(\lambda)]'$, $\mathbf{W} = \{w_{db}\}$, and $\mathbf{r} = [r_1, \dots, r_D]'$. The symbols \odot and \oslash denote component-wise multiplication and division respectively.

The complete-data space \mathbf{X}_{IA} for ML-IA is

$$\mathbf{X}_{IA} = \{ \{N_{db}\}, \{R_d\}, \{N_{db}^o\} \}, \quad \begin{matrix} d = 1, \dots, D \\ b = 1, \dots, B \end{matrix}$$

which includes not only the components of \mathbf{X}_{IB} , but also the attenuated counts N_{db}^o . These have independent Poisson distributions with mean $(1 - \alpha_d)w_{db}\lambda_b$. Thus,

$$\log f_{IA}(\mathbf{X}_{IA}; \bar{\lambda}) = \log f_{IB}(\mathbf{X}_{IB}; \bar{\lambda}) + \log f_o(\mathbf{N}^o; \bar{\lambda}), \quad (10)$$

where

$$\log f_o(\mathbf{N}^o; \bar{\lambda}) =$$

$$\sum_{d=1}^D \sum_{b=1}^B \left(-(1 - \alpha_d)w_{db}\bar{\lambda}_b + N_{db}^o \log((1 - \alpha_d)w_{db}\bar{\lambda}_b) \right). \quad (11)$$

Since N_{db}^o makes *no contribution* to the measurements \mathbf{Y} ,

$$E\{N_{db}^o | \mathbf{y}; \lambda\} = (1 - \alpha_d)w_{db}\lambda_b. \quad (12)$$

Combining equations (10)-(12) and (1) yields

$$Q_{IA}(\bar{\lambda}; \lambda) = Q_{IB}(\bar{\lambda}; \lambda) + \sum_{d=1}^D \sum_{b=1}^B \left(-(1 - \alpha_d)w_{db}\bar{\lambda}_b + (1 - \alpha_d)w_{db}\lambda_b \log((1 - \alpha_d)w_{db}\bar{\lambda}_b) \right).$$

Maximizing $Q_{IA}(\bar{\lambda}; \lambda)$ over $\bar{\lambda}$ yields the ML-IA iteration $\lambda^{i+1} = G_{IA}(\lambda^i)$, where

$$G_{IA}(\lambda) \triangleq \lambda - \lambda \odot (\mathbf{W}'\alpha) \odot (\mathbf{W}'\mathbf{1}) + \lambda \odot [\mathbf{W}'(\alpha \odot \mathbf{y} \odot \hat{\mathbf{y}}(\lambda))] \odot [\mathbf{W}'\mathbf{1}].$$

The vector $\mathbf{1}$ is the $D \times 1$ vector of 1's.

Note that the fixed point(s) of G_{IA} and G_{IB} are identical, and if there were no attenuation, i.e. $\alpha = \mathbf{1}$, then the two algorithms would be identical.

B. Convergence rates

Since the complete-data space \mathbf{X}_{IA} contains the attenuated events N_{db}^o that do not contribute to the measurements, conditions (ii)-(iv) of Theorem 2 are satisfied. Before invoking Theorem 2 to conclude that ML-IB is faster, we must verify condition (i) of Theorem 2, including strict concavity, convergence of the algorithms, and convergence to an interior point. Under the assumption of strict concavity, discussed further below, one can apply the same arguments as in the appendix of [4] to show that ML-IA and ML-IB converge globally to the same estimate. (Strict concavity is a sufficient condition for convergence, but it may not be necessary.) From (4) and (8),

$$L(\lambda) = \log g(\mathbf{y}; \lambda) = \sum_{d=1}^D (-\hat{y}_d(\lambda) + y_d \log(\hat{y}_d(\lambda))), \quad (13)$$

so the Hessian of the likelihood is

$$\nabla^2 L(\lambda) = -\mathbf{W}' \text{diag}\{\mathbf{y} \odot \hat{\mathbf{y}}^2(\lambda) \odot \alpha^2\} \mathbf{W}. \quad (14)$$

For strict concavity it is sufficient to have $\alpha_d y_d > 0, \forall d$, provided \mathbf{W} has full column rank. In the presence of accidental coincidences, all PET detectors record nonzero coincidences with very high probability. For an appropriate sampling scheme and a well-designed PET system, \mathbf{W} should have full column rank. If \mathbf{W} is not full rank, such as when "too many" pixels are used, then the ML estimate is not unique, and the likelihood criterion is inappropriate. In such cases a penalized likelihood estimate is preferable, as we discuss in Section V. We conclude then, that if the \mathbf{W} is full rank, and if the ML estimate is strictly positive, then the ML-IB algorithm converges faster than ML-IA.

Unfortunately, in most cases the ML estimate in PET will have components that are zero [13], i.e., not in the interior of \mathcal{D} . Strictly speaking, the above analysis is inconclusive for such examples. How likely is it that ML-IB, with its smaller complete-data space, would converge slower than ML-IA simply because some of the components converge to zero? In the next section we explore this question by considering a one-dimensional analogue.

C. Scalar example

One can obtain some insight into the convergence behavior of these two algorithms by considering the following scalar version of the problem. Suppose the measurement model is:

$$y \sim \text{Poisson}(a\lambda + r)$$

where the attenuation $a \in (0, 1)$ and the accidental coincidence rate $r \geq 0$ are known. In this case, the ML estimate for emission rate λ over $\mathcal{D} = \{\lambda : \lambda \geq 0\}$ is given by:

$$\hat{\lambda} = \max \left\{ 0, \frac{y - r}{a} \right\},$$

a truncated subtraction. The ML-IA and ML-IB algorithms are given respectively by the maps

$$G_{IA}(\lambda) = (1 - a)\lambda + \lambda a \frac{y}{a\lambda + r}$$

and

$$G_{IB}(\lambda) = \lambda \frac{y}{a\lambda + r}.$$

Note that in the absence of attenuation ($a = 1$), the two algorithms are identical. One can also verify that both algorithms are globally convergent if $\lambda^0 > 0$. Differentiating:

$$\begin{aligned} \frac{d}{d\lambda} G_{IA}(\lambda) &= (1 - a) + \frac{ayr}{(a\lambda + r)^2} \\ \frac{d}{d\lambda} G_{IB}(\lambda) &= \frac{yr}{(a\lambda + r)^2}, \end{aligned}$$

so in particular

$$\begin{aligned} \rho_B &= \frac{d}{d\lambda} G_{IB}(\lambda)|_{\hat{\lambda}} = \min \left\{ \frac{r}{y}, \frac{y}{r} \right\} \\ \rho_A &= \frac{d}{d\lambda} G_{IA}(\lambda)|_{\hat{\lambda}} = 1 - a + a \frac{d}{d\lambda} G_{IB}(\lambda)|_{\hat{\lambda}} \\ &= 1 + a(\rho_B - 1) \geq \rho_B, \end{aligned} \quad (15)$$

showing that the root-convergence factor for ML-IB is smaller than that of ML-IA. Does ML-IB converge faster? There are three cases to consider.

Case 1: If $y > r$, then both estimates converge to $\hat{\lambda} > 0$, at asymptotic rates governed by the Linear Convergence Theorem, so by (15), ML-IB converges faster.

Case 2: If $y \leq r$, then both estimates converge to $\hat{\lambda} = 0$, on the boundary of \mathcal{D} , so at first it seems that the Linear Convergence Theorem does not apply. However, if $r > 0$ then we can actually make the object domain slightly larger, say: $\mathcal{D}_- = \{\lambda : \lambda \geq -\frac{1}{2}r/a\}$, since G_{IA} and G_{IB} are both differentiable on \mathcal{D}_- . Directly applying³ the Linear Convergence Theorem to G_{IA} and G_{IB} using (15) shows that if $y < r$, then ML-IB converges faster than ML-IA even though the ML estimate is 0!

Case 3: If $y = r$, then $\rho_A = \rho_B = 1$, so the asymptotic convergence rate is not well defined by the Linear Convergence Theorem. However, since y is an integer number of counts, and r is a real number, the outcome $y = r$ seems rather unlikely in practice. For a non-asymptotic comparison, one can verify that if $\lambda > 0$, then

$$|G_{IB}(\lambda) - \hat{\lambda}| \leq |G_{IA}(\lambda) - \hat{\lambda}|, \quad (16)$$

so the ML-IB algorithm takes larger steps towards the ML estimate than the ML-IA algorithm. Therefore, even though the convergence is sub-linear when $y = r$, the ML-IB algorithm will converge faster in a sub-linear sense.

In summary, we have shown that under this scalar model, ML-IB usually has faster asymptotic convergence rate than ML-IA, and always takes larger steps (16). It is

³We cannot apply Theorem 1 to this larger domain since Q_{IA} and Q_{IB} are not differentiable at 0.

difficult to predict what the analogous boundary situations would be in higher dimensions. The fact that there exists a situation where $\rho(\nabla G) = 1$ even in the scalar case suggests that a comprehensive rigorous comparison of ML-IB and ML-IA will be difficult to obtain.

IV. JOINT ESTIMATION OF α AND λ

As in [1], the above discussion assumed that the survival probabilities α_d were known exactly. In practice, one estimates these survival probabilities from transmission measurements acquired prior to the emission scan. Let M_d denote the transmission measurement, with realization m_d , for the d th detector pair. It is reasonable to assume that the M_d 's have independent Poisson distributions with mean $t_d \alpha_d$, where t_d is proportional to the transmission scan time⁴ and the efficiency of the d th detector pair. The t_d factors are determined by a "blank scan," i.e. a transmission scan without the patient in the scanner. The conventional approach is simply to estimate α_d by m_d/t_d , and to "precorrect" the emission measurements. The approach suggested in [8] is to jointly estimate λ and α from \mathbf{y} and $\mathbf{m} = [m_1, \dots, m_D]'$ by maximizing the joint log-likelihood:

$$L(\theta) = \log g(\mathbf{y}, \mathbf{m}; \lambda, \alpha) = \log g(\mathbf{y}; \lambda, \alpha) + \log g(\mathbf{m}; \alpha)$$

$$= \sum_{d=1}^D (-\hat{y}_d(\theta) + y_d \log(\hat{y}_d(\theta)) - t_d \alpha_d + m_d \log(t_d \alpha_d)),$$

where $\theta = \{\lambda, \alpha\}$ and $\hat{y}_d(\theta) = \sum_{b=1}^B \alpha_d w_{db} \lambda_b + r_d$. It follows that

$$-\nabla^2 L(\theta) =$$

$$\begin{bmatrix} \mathbf{W}' & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \left(\mathbf{K}(\theta) + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \text{diag}\{\mathbf{m} \odot \alpha^2\} \end{bmatrix} \right) \begin{bmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix},$$

where

$$\mathbf{K}(\theta) =$$

$$\begin{bmatrix} \text{diag}\{\mathbf{y} \odot \hat{\mathbf{y}}^2(\theta) \odot \alpha^2\} & \text{diag}\{\mathbf{1} - \mathbf{y} \odot \mathbf{r} \odot \hat{\mathbf{y}}^2(\theta)\} \\ \text{diag}\{\mathbf{1} - \mathbf{y} \odot \mathbf{r} \odot \hat{\mathbf{y}}^2(\theta)\} & \text{diag}\{\mathbf{y} \odot \hat{\mathbf{y}}^2(\theta) \odot (\mathbf{W}\lambda)^2\} \end{bmatrix}.$$

Thus, $L(\theta)$ is not necessarily strictly concave, and convergence of the two EM algorithms discussed below remains an open problem.

As in Section III, we again have the option of including or excluding the attenuated emissions N_{db}^o from the complete-data space, leading to two algorithms we denote ML-JA and ML-JB.

For ML-JB, the complete data consists of

$$\mathbf{X}_{JB} = \{ \{N_{db}\}, \{R_d\}, \{M_d\} \}, \quad d = 1, \dots, D, \quad b = 1, \dots, B,$$

⁴Here, we ignore the statistical uncertainty in the accidental coincidences measurement, the statistical uncertainty in the blank scan, and the contribution of accidental coincidences to the transmission measurements. For a more complete treatment, see [8].

so

$$\begin{aligned} \log f_{\text{JB}}(\mathbf{X}_{\text{JB}}; \bar{\theta}) = & \sum_{d=1}^D \sum_{b=1}^B (-\bar{\alpha}_d w_{db} \bar{\lambda}_b + N_{db} \log(\bar{\alpha}_d w_{db} \bar{\lambda}_b)) \\ & + \sum_{d=1}^D (-r_d + R_d \log(r_d)) + \sum_{d=1}^D (-t_d \bar{\alpha}_d + M_d \log(t_d \bar{\alpha}_d)). \end{aligned}$$

It is easily verified that

$$\begin{aligned} E\{N_{db} | \mathbf{y}, \mathbf{m}; \theta\} &= y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\theta)}, \\ E\{R_d | \mathbf{y}, \mathbf{m}; \theta\} &= y_d \frac{r_d}{\hat{y}_d(\theta)}, \\ E\{M_d | \mathbf{y}, \mathbf{m}; \theta\} &= m_d; \end{aligned}$$

therefore,

$$\begin{aligned} Q_{\text{JB}}(\bar{\theta}; \theta) = & \sum_{d=1}^D \sum_{b=1}^B \left(-\bar{\alpha}_d w_{db} \bar{\lambda}_b + y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\theta)} \log(\bar{\alpha}_d w_{db} \bar{\lambda}_b) \right) \\ & + \sum_{d=1}^D \left(-r_d + y_d \frac{r_d}{\hat{y}_d(\theta)} \log(r_d) \right) + \sum_{d=1}^D (-t_d \bar{\alpha}_d + m_d \log(t_d \bar{\alpha}_d)) \end{aligned}$$

Setting the derivatives of $Q_{\text{JB}}(\cdot, \theta)$ to zero yields the following equations for the M-step:

$$\begin{aligned} 0 &= \sum_{d=1}^D \left(-\bar{\alpha}_d w_{db} + y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\theta)} / \bar{\lambda}_b \right), \\ 0 &= \sum_{b=1}^B \left(-w_{db} \bar{\lambda}_b + y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\theta)} / \bar{\alpha}_d \right) - t_d + m_d / \bar{\alpha}_d, \end{aligned}$$

yielding the pseudo-iteration:

$$\begin{aligned} \lambda^{i+1} &= \lambda^i \odot [\mathbf{W}'(\alpha^i \odot \mathbf{y} \odot \hat{\mathbf{y}}(\theta^i))] \odot [\mathbf{W}'\alpha^{i+1}], \\ \alpha^{i+1} &= [\mathbf{y} \odot \hat{\mathbf{y}}(\theta^i) \odot \alpha^i \odot (\mathbf{W}\lambda^i) + \mathbf{m}] \odot [\mathbf{W}\lambda^{i+1} + \mathbf{t}], \end{aligned}$$

where $\mathbf{t} = [t_1, \dots, t_D]'$. This set of equations is coupled in λ^{i+1} and α^{i+1} , and no analytical solution seems likely. However, they can form the basis for a GEM algorithm [8, 14].

Fortunately, the equations become uncoupled when using the less parsimonious complete-data space for ML-JA. Let

$$\mathbf{X}_{\text{JA}} = \{ \{N_{db}\}, \{R_d\}, \{M_d\}, \{N_{db}^o\} \}, \quad d = 1, \dots, D, \quad b = 1, \dots, B,$$

then

$$\begin{aligned} \log f_{\text{JA}}(\mathbf{X}_{\text{JA}}; \bar{\theta}) &= \log f_{\text{JB}}(\mathbf{X}_{\text{JB}}; \bar{\theta}) \\ &+ \sum_{d=1}^D \sum_{b=1}^B (- (1 - \bar{\alpha}_d) w_{db} \bar{\lambda}_b + N_{db}^o \log((1 - \bar{\alpha}_d) w_{db} \bar{\lambda}_b)). \end{aligned}$$

Since N_{db}^o is independent of the measurements,

$$E\{N_{db}^o | \mathbf{y}, \mathbf{m}; \theta\} = (1 - \alpha_d) w_{db} \lambda_b,$$

thus

$$\begin{aligned} Q_{\text{JA}}(\bar{\theta}; \theta) &= Q_{\text{JB}}(\bar{\theta}; \theta) + \sum_{d=1}^D \sum_{b=1}^B (- (1 - \bar{\alpha}_d) w_{db} \bar{\lambda}_b \\ &+ (1 - \alpha_d) w_{db} \lambda_b \log((1 - \bar{\alpha}_d) w_{db} \bar{\lambda}_b)). \end{aligned}$$

Setting the derivatives of $Q_{\text{JA}}(\cdot, \theta)$ to zero yields the following equations for the M-step:

$$\begin{aligned} 0 &= \sum_{d=1}^D \left(-\bar{\alpha}_d w_{db} + y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\theta)} / \bar{\lambda}_b \right) \\ &+ \sum_{d=1}^D (- (1 - \bar{\alpha}_d) w_{db} + (1 - \alpha_d) w_{db} \lambda_b / \bar{\lambda}_b), \\ 0 &= \sum_{b=1}^B \left(-w_{db} \bar{\lambda}_b + y_d \frac{\alpha_d w_{db} \lambda_b}{\hat{y}_d(\theta)} / \bar{\alpha}_d \right) \\ &+ (-t_d + m_d / \bar{\alpha}_d) \\ &+ \sum_{b=1}^B (w_{db} \bar{\lambda}_b + (1 - \alpha_d) w_{db} \lambda_b / (1 - \bar{\alpha}_d)). \end{aligned}$$

This set of equations is uncoupled in λ^{i+1} and α^{i+1} , and yields the ML-JA iteration:

$$\begin{aligned} \lambda^{i+1} &= \lambda^i - \lambda^i \odot (\mathbf{W}'\alpha^i) \odot (\mathbf{W}'\mathbf{1}) \\ &+ \lambda^i \odot [\mathbf{W}'(\alpha^i \odot \mathbf{y} \odot \hat{\mathbf{y}}(\theta^i))] \odot [\mathbf{W}'\mathbf{1}], \end{aligned}$$

and

$$\begin{aligned} \mathbf{t} &= [(\mathbf{W}\lambda^i) \odot \alpha^i \odot \mathbf{y} \odot \hat{\mathbf{y}}(\theta^i) + \mathbf{m}] \odot \alpha^{i+1} \\ &- [(\mathbf{W}\lambda^i) \odot (\mathbf{1} - \alpha^i)] \odot (\mathbf{1} - \alpha^{i+1}), \end{aligned}$$

where the latter is an easily solved quadratic in α^{i+1} . The resulting joint emission/transmission estimation algorithm is only slightly more computationally expensive per iteration than the ML-IA or ML-IB algorithms of [1], yet it accounts for the statistical uncertainty in both the emission measurements and the transmission measurements, unlike ML-IA and ML-IB.

V. DISCUSSION

We have shown that smaller complete-data spaces yield EM algorithms with faster asymptotic convergence. This theoretical result, combined with the empirical results in [1] suggests strongly that the ML-IB algorithm should be used in practice over the ML-IA algorithm. The heuristic explanation for this is that the complete-data space for ML-IA includes the attenuated events that make no contribution to the measurements. Since EM algorithms are notorious for slow convergence, this comparison has practical importance. Even a small decrease in the root-convergence factor can significantly reduce the required number of iterations.

We have also shown that the story gets more complicated if one wants to jointly estimate both the emission

and the transmission parameters. In this case, although theoretically ML-JB would converge faster than ML-JA, the M-step of ML-JB seems intractable. We are currently investigating a space-alternating approach that may circumvent this problem [14]. Meanwhile, it appears that the best strategy is the following: *use the smallest complete-data space that results in a tractable maximization step.*

Several investigators have shown that more appealing images are produced by regularizing the ML estimate by including a penalty term or Bayesian "prior" [12, 15–18]. In principle, our Theorems 1 and 2 directly generalize to the case where concave penalties such as those discussed in [18] are added to the likelihood, again supporting the conclusion that smaller complete-data spaces correspond to faster convergence. There is one important caveat however: except in the trivial case of independent priors, the maximization steps of penalized EM algorithms become intractable due to the coupling introduced by the penalties. Consequently, the algorithms for the penalized case are usually of the generalized EM (GEM) type [3, 16]. GEM algorithms only provide an increase in $Q(\theta, \theta^i)$ at each iteration, rather than truly maximizing Q . Therefore, GEM algorithms do not usually satisfy condition (ii) of our Theorem 1. They are also usually not globally convergent unless line-searches are employed [18]. These factors inhibit making formal statements about asymptotic convergence rates for penalized likelihood algorithms. We have implemented penalized-likelihood algorithms based on Hebert's GEM strategy [16] for both the ML-IA and ML-IB complete-data spaces. We have also implemented both ML-IA and ML-IB with sieve constraints [1, 19, 8]. We found empirically that the penalized ML-IB algorithm converged substantially more rapidly, in terms of both likelihood increase and apparent image contrast. These empirical results are further motivation for using smaller complete-data spaces where possible.

VI. ACKNOWLEDGEMENT

The authors gratefully acknowledge discussions with P. Chiao and A. Hero.

REFERENCES

- [1] D. G. Politte and D. L. Snyder. Corrections for accidental coincidences and attenuation in maximum-likelihood image reconstruction for positron-emission tomography. *IEEE Transactions on Medical Imaging*, 10(1):82–89, March 1991.
- [2] L. A. Shepp and Y. Vardi. Maximum likelihood reconstruction for emission tomography. *IEEE Transactions on Medical Imaging*, 1(2):113–122, October 1982.
- [3] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society Series B*, 39(1):1–38, 1977.
- [4] K. Lange and R. Carson. EM reconstruction algorithms for emission and transmission tomography. *Journal of Computer Assisted Tomography*, 8(2):306–316, April 1984.
- [5] M. I. Miller, D. L. Snyder, and T. R. Miller. Maximum-likelihood reconstruction for single-photon emission computed-tomography. *IEEE Transactions on Nuclear Science*, 32(1):769–778, February 1985.
- [6] N. H. Clinthorne, X. H. Wang, and J. A. Fessler. Multi-energy maximum-likelihood reconstruction algorithms for SPECT and PET. *Journal of Nuclear Medicine (Abstract Book)*, 33(5):831, 1992.
- [7] M. E. Daube-Witherspoon, R. E. Carson, Y. Yan, and T. K. Yap. Scatter correction in maximum likelihood reconstruction of PET data. In *Abstract Book of the 1992 IEEE Nuclear Science Symposium and Medical Imaging Conference*, 1992.
- [8] N. H. Clinthorne, J. A. Fessler, G. D. Hutchins, and W. L. Rogers. Joint maximum likelihood estimation of emission and attenuation densities in PET. In *Conference Record of the 1991 IEEE Nuclear Science Symposium and Medical Imaging Conference*, volume 3, pages 1927–1932, 1991.
- [9] M. Segal and E. Weinstein. The cascade EM algorithm. *Proceedings of the IEEE*, 76(10):1388–1390, October 1988.
- [10] J. M. Ortega and W. C. Rheinboldt. *Iterative Solution of Nonlinear Equations in Several Variables*. Academic Press, New York, 1970.
- [11] E. L. Lehmann. *Theory of Point Estimation*. Wiley, New York, 1983.
- [12] P. J. Green. On use of the EM algorithm for penalized likelihood estimation. *Journal of the Royal Statistical Society Series B*, 52(3):443–452, 1990.
- [13] C. L. Byrne. Iterative image reconstruction reconstruction algorithms based on cross-entropy minimization, 1991. Submitted to IEEE Trans. Image Proc.
- [14] J. A. Fessler and A. O. Hero. Complete-data spaces and generalized em algorithms. In *Proc. IEEE Conf. on Acoustics, Speech, and Signal Processing*, 1993. To appear.
- [15] S. Geman and D. E. McClure. Bayesian image analysis: An application to single photon emission tomography. *Proc. Amer. Stat. Ass. Stat. Comp. Sect.*, pages 12–18, 1985.
- [16] T. Hebert and R. Leahy. A generalized EM algorithm for 3-D Bayesian reconstruction from Poisson data using Gibbs priors. *IEEE Transactions on Medical Imaging*, 8(2):194–202, June 1989.
- [17] C. T. Chen, V. E. Johnson, W. H. Wong, X. Hu, and C. E. Metz. Bayesian image reconstruction in positron emission tomography. *IEEE Transactions on Nuclear Science*, 37(2):636–641, April 1990.
- [18] K. Lange. Convergence of EM image reconstruction algorithms with Gibbs smoothing. *IEEE Transactions on Medical Imaging*, 9(4):439–446, December 1990. Corrections, June 1991 TMI.
- [19] D. L. Snyder, M. I. Miller, L. J. Thomas, and D. G. Politte. Noise and edge artifacts in maximum-likelihood reconstructions for emission tomography. *IEEE Transactions on Medical Imaging*, 6(3):228–238, September 1987.