# A Cheaper Way to Compute Generalized Cross-Validation as a Stopping Rule for Linear Stationary Iterative Methods

## Reginaldo J. SANTOS and Álvaro R. DE PIERRO

We apply generalized cross-validation (GCV) as a stopping rule for general linear stationary iterative methods for solving very large-scale, ill-conditioned problems. We present a new general formula for the influence operator for these methods and, using this formula and a Monte Carlo approach, we show how to compute the GCV function at a cheaper cost. Then we apply our approach to a well known iterative method (ART) with simulated data in positron emission tomography (PET).

**Key Words:** Emission tomography; Ill-posed problems; Parameter estimation.

## 1. INTRODUCTION

We consider the problem of estimating a solution $\mathbf{x}$ of

$$\mathbf{Ax} + \boldsymbol{\epsilon} = \mathbf{b}, \tag{1.1}$$

where $\mathbf{A}$ is a real $m \times n$ matrix, $\mathbf{b}$ is an $m$-vector of observations and $\boldsymbol{\epsilon}$ is a random vector with

$$E\boldsymbol{\epsilon} = 0, \quad E\boldsymbol{\epsilon}\boldsymbol{\epsilon}^t = \sigma^2 \mathbf{I}, \tag{1.2}$$

where $E$ denotes the expectation and $\mathbf{I}$ is the identity matrix and the, possibly unknown, variance $\sigma^2$.

If $\mathbf{A}$ is a full-rank matrix, the Gauss–Markov theorem states that the least squares estimator of (1.1), that is, $\tilde{\mathbf{x}} = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t \mathbf{b}$, is the best unbiased linear estimator of $\mathbf{x}$, meaning that it is the minimum variance estimator (see, e.g., Silvey 1970). However, if $\mathbf{A}$

Reginaldo J. Santos is Associate Professor, Department of Mathematics, Federal University of Minas Gerais, CP 702, 30123-970, Belo Horizonte, MG, Brazil (E-mail: regi@mat.ufmg.br). Álvaro R. De Pierro is Professor, Department of Applied Mathematics, State University of Campinas, CP 6065, 13081-970, Campinas, SP, Brazil (E-mail: aldep@math.ucla.edu).

is ill-conditioned, this minimum variance is still large. It is well known that, if we allow the estimator to be biased, the variance could be drastically reduced (see van der Sluis and van der Vorst 1987; 1990).

One way to do this is by considering solutions of regularized problems of the form

$$\text{minimize} \quad ||\mathbf{Ax} - \mathbf{b}||^2 + \lambda \mathbf{x}^t \mathbf{Bx} , \tag{1.3}$$

for a positive real number $\lambda$, where $\mathbf{B}$ is an $n \times n$ matrix that introduces a priori information on the problem. For example, if we take $\mathbf{B} = \mathbf{L}^t\mathbf{L}$, where $\mathbf{L}$ is the discretization of first-order derivatives, increasing $\lambda$ restricts first-order variations of $x$.

Alternatively, one can apply an iterative method that is convergent to

$$\mathbf{MAx} = \mathbf{Mb},$$

where $\mathbf{M}$ can be, for example, equal to $\mathbf{A}^t$. In this case, a regularization effect is obtained by choosing as a "solution" an early iterate of the method. Fleming (1990) and Santos (1996) proved that, for stationary linear iterative methods this approach is equivalent, in some sense, to (1.3). Now, the role of the parameter $\lambda$ is played by the iteration number.

A crucial question is how to choose the regularization parameter $\lambda$ in (1.3) or the iterate $k$ in the case of iterative methods, in such a way that the "information loss" is minimized. One possibility is to approximate a value of $\lambda$ such that the average of the mean square error when estimating $\mathbf{Ax}$ is minimum, that is, $\lambda$ solves the problem

$$\text{minimize} \quad ET(\lambda) , \tag{1.4}$$

where

$$T(\lambda) = \frac{1}{m}||\mathbf{Ax} - \mathbf{Ax}_\lambda||^2 , \tag{1.5}$$

and $\mathbf{x}_\lambda$ and $\mathbf{x}$ are solutions of (1.3) and (1.1), respectively. In spite of the fact that we do not know the values of $\mathbf{x}$ and $\sigma^2$, there exist methods that provide good approximations for the solution of (1.4). Probably, the most popular of these methods is generalized cross-validation (GCV).

In the case of Equation (1.1) and if $x_\lambda$ is a solution of (1.3), defined by

$$V(\lambda) = \frac{\frac{1}{m}||\mathbf{b} - \mathbf{Ax}_\lambda||^2}{[\frac{1}{m}Tr(\mathbf{I} - \mathbf{A}(\lambda))]^2} \tag{1.6}$$

(Craven and Wahba 1979; Golub, Heath, and Wahba 1979), where $\mathbf{A}(\lambda)$ is the influence operator defined as

$$\mathbf{A}(\lambda)\mathbf{b} = \mathbf{Ax}_\lambda. \tag{1.7}$$

GCV chooses the regularization parameter by minimizing $V(\lambda)$. This estimate is a rotationally invariant version of Allen's PRESS, or cross-validation (Allen 1974). Craven

and Wahba (1979) and Golub, Heath, and Wahba (1979) proved that, when $x_\lambda$ is the solution of (1.3), if $\lambda_0$ is a minimizer of $ET(\lambda)$ and $\tilde{\lambda}$ the minimizer of $EV(\lambda)$, then

$$\frac{ET(\tilde{\lambda})}{ET(\lambda_0)} \leq \frac{1 + h(\lambda_0)}{1 - h(\tilde{\lambda})},\tag{1.8}$$

where

$$h(\lambda) = \left(2\mu_1(\lambda) + \frac{\mu_1(\lambda)^2}{\mu_2(\lambda)}\right)\frac{1}{(1 - \mu_1(\lambda))^2},$$

$$\mu_1(\lambda) = (1/m)Tr\mathbf{A}(\lambda),$$

and

$$\mu_2(\lambda) = (1/m)Tr\mathbf{A}^2(\lambda).$$

Therefore, if $h(\lambda_0)$ and $h(\tilde{\lambda})$ are small enough, the mean square error for $\tilde{\lambda}$ is not much greater than the mean square error for $\lambda_0$. Also, Golub, Heath, and Wahba (1979) proved that, under reasonable assumptions,

$$\frac{ET(\tilde{\lambda})}{ET(\lambda_0)} \downarrow 1,$$

when $n$ tends to infinity. Golub and von Matt (1997) proposed an iterative method to approximate the GCV function for large scale problems.

Wahba (1987) showed how to apply GCV as a stopping rule for linear stationary iterative methods of the form $\mathbf{x}^{k+1} = \mathbf{x}^k + \tilde{\mathbf{M}}\mathbf{A}^t(\mathbf{b} - \mathbf{A}\mathbf{x}^k)$, $k = 0, 1, 2, \ldots$, if $\tilde{\mathbf{M}}$ is symmetric and positive definite and the starting point is $\mathbf{x}^0 = 0$. She used the SVD of $\tilde{\mathbf{M}}^{1/2}\mathbf{A}$.

GCV as stopping rule of iterative methods appeared repeatedly in the engineering literature (see, e.g., Reeves 1992, 1994; Perry and Reeves 1994). However, the nonlinearity of the influence operator has not been taken into account.

Our main goal in this article is the application of GCV as a stopping rule for general linear stationary iterative methods to large scale ill-conditioned problems like positron emission tomography (PET); that is our reference problem. The algorithm we chose to illustrate our approach is ART [from algebraic reconstruction technique; Herman (1980)], a very well-known iterative algorithm dating back to Kaczmarz (1937). Two problems have to be solved in order to apply GCV to ART in PET. The existing results in the GCV literature are valid when the influence matrix is symmetric positive semidefinite and this is not the case for ART. On the other hand, if the starting point is not zero, the influence operator for ART is an affine transformation. The second problem is how to compute $Tr(\mathbf{I} - \mathbf{A}(\lambda))$ at a reasonable cost.

Section 2 shows that an inequality like (1.8) is still valid when $\mathbf{A}(\lambda)$ is affine. Section 3 gives a new general formula for computing the influence operator for convergent linear stationary iterative methods. We show how to apply this formula to approximate the trace of the linear part of the influence operator saving the cost equivalent to a product of $\mathbf{A}$ by a

vector in the computation of the GCV function proposed in Section 2. Section 4 proves some properties of a Monte Carlo method used to approximate the trace of a general matrix; these proofs are slight but useful extensions of those given by Girard (1989). Section 5 briefly describes the PET problem and the implementation of ART and we present simulations and the numerical results. Section 6 presents some concluding remarks.

## 2. GCV FOR AFFINE INFLUENCE OPERATORS

Let $\{\mathbf{x}_\lambda\}$ be a general family of estimates of the solution of (1.1). For example, this family can be generated by the iterates of an iterative convergent method. The point of departure for GCV could be to find a good estimate for the mean square error of $\mathbf{Ax}$, that is, (1.5). In order to approximate $T(\lambda)$, we have to start with the mean square residual

$$U(\lambda) = \frac{1}{m}||\mathbf{b} - \mathbf{Ax}_\lambda||^2. \tag{2.1}$$

Throughout this article we will assume that the influence operator, defined by $\mathbf{A}(\lambda)\mathbf{b} = \mathbf{Ax}_\lambda$, has the form

$$\mathbf{A}(\lambda)(\mathbf{b}) = \mathbf{A}_0(\lambda)\mathbf{b} + \mathbf{b}_0(\lambda), \tag{2.2}$$

where $\mathbf{A}_0(\lambda)$ is an $m \times m$ matrix and $\mathbf{b}_0(\lambda)$ is an $m$-vector.

Let us now define

$$\tilde{\mu}_1(\lambda) = \frac{1}{m}Tr(\mathbf{A}_0(\lambda)), \tag{2.3}$$

and

$$\tilde{\mu}_2(\lambda) = \frac{1}{m}Tr(\mathbf{A}_0(\lambda)^t\mathbf{A}_0(\lambda)). \tag{2.4}$$

Now, for our problem the GCV function is defined as

$$V(\lambda) = \frac{\frac{1}{m}||\mathbf{b} - \mathbf{Ax}_\lambda||^2}{[\frac{1}{m}Tr(\mathbf{I} - \mathbf{A}_0(\lambda))]^2} = \frac{U(\lambda)}{(1 - \tilde{\mu}_1(\lambda))^2}. \tag{2.5}$$

**Theorem 1.** *Let $\{\mathbf{x}_\lambda\}$ be a family of estimates of the solution of (1.1), with $\epsilon$ satisfying (1.2). If $\mathbf{A}(\lambda)$ has the form (2.2), $V(\lambda)$ is the GCV function defined by (2.5), $\lambda_0$ and $\lambda_1$ are global minima of $ET(\lambda)$ and $EV(\lambda)$, respectively, then*

$$\frac{ET(\lambda_1)}{ET(\lambda_0)} \leq \frac{1 + \tilde{h}(\lambda_0)}{1 - \tilde{h}(\lambda_1)}, \tag{2.6}$$

*where*

$$\tilde{h}(\lambda) = \frac{1}{(1 - \tilde{\mu}_1(\lambda))^2}\left(\frac{\tilde{\mu}_1(\lambda)^2}{\tilde{\mu}_2(\lambda)} + 2|\tilde{\mu}_1(\lambda)| + \tilde{\mu}_1(\lambda)^2\right),$$

*$\tilde{\mu}_1(\lambda)$, and $\tilde{\mu}_2(\lambda)$ are defined by (2.3) and (2.4), respectively.*

**Proof:** Similar to Theorem 1 in Golub, Heath, and Wahba (1979).  □

## 3. RANDOMIZED GCV AS A STOPPING RULE FOR LSIM

Our main goal in this section is to determine a stopping rule for any convergent linear stationary iterative method (LSIM) for solving (1.1). Now, each iterate $x^k$ is an estimate of the solution of (1.1) and $k$ plays the role of $\lambda$. Section 2 shows that GCV is applicable to affine influence operators and this is the first step towards using GCV for an LSIM.

One of the difficulties in applying GCV to large-scale problems is the calculation of the trace in the denominator of (2.5). The idea of a randomized GCV is to replace the trace in the denominator of (2.5) by an estimate of the trace as was proposed by Girard (1989) for the case of Tikhonov's regularization. In that case the influence operator is linear and symmetric.

Applying Girard's idea in our context we have

$$\frac{1}{m}Tr(\mathbf{I} - \mathbf{A}_0(k)) \approx \frac{\mathbf{w}^t\mathbf{w} - \mathbf{w}^t\mathbf{A}_0(k)\mathbf{w}}{\mathbf{w}^t\mathbf{w}},$$

where $\mathbf{w} = (w_1, \ldots, w_m)^t \in \mathbb{R}^m$ is a random vector from a normal distribution with standard deviation 1, that is, $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{n \times n})$. $\mathbf{A}_0(k)\mathbf{w}$ is just the matrix $\mathbf{A}$ applied to the $k$th iteration of (3.3) using $\mathbf{b} = \mathbf{w}$ and zero as a starting point. However, using Theorem 2 below, we obtain a cheaper way to compute the approximation of $Tr(\mathbf{I} - \mathbf{A}_0(k))$ in the randomized GCV function.

Consider a convergent LSIM of the form

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \tilde{\mathbf{M}}\mathbf{A}^t(\mathbf{b} - \mathbf{A}x^k), \qquad k = 0, 1, 2, \ldots. \tag{3.1}$$

If $\tilde{\mathbf{M}}$ is symmetric and positive definite, Wahba (1987) proved that if the starting point is $\mathbf{x}^0 = 0$ the influence matrix is given by

$$\mathbf{A}(k) = \mathbf{I} - (\mathbf{I} - \mathbf{A}\tilde{\mathbf{M}}\mathbf{A}^t)^k. \tag{3.2}$$

The following theorem generalizes (3.2) to any LSIM.

**Theorem 2.** *Consider*

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \mathbf{M}(\mathbf{b} - \mathbf{A}\mathbf{x}^k), \qquad k = 0, 1, 2, \ldots, \tag{3.3}$$

*where $\mathbf{M}$ is an $n \times m$ matrix, such that (3.3) is a convergent method. Then, the influence operator $\mathbf{A}(k)$, for this method is given by*

$$\mathbf{A}(k)(\mathbf{b}) = [\mathbf{I} - (\mathbf{I} - \mathbf{A}\mathbf{M})^k]\mathbf{b} + \mathbf{A}\mathbf{x}_1^0 + \mathbf{A}(\mathbf{I} - \mathbf{M}\mathbf{A})^k\mathbf{x}_2^0, \tag{3.4}$$

*where $\mathbf{x}_1^0 \in \mathcal{N}(\mathbf{M}\mathbf{A})$ and $\mathbf{x}_2^0 \in \mathcal{R}(\mathbf{M}\mathbf{A})$ are such that $\mathbf{x}^0 = \mathbf{x}_1^0 + \mathbf{x}_2^0$. $\mathcal{N}(\mathbf{M}\mathbf{A})$ and $\mathcal{R}(\mathbf{M}\mathbf{A})$ denote the null space and the range of $\mathbf{M}\mathbf{A}$, respectively.*

**Proof:** By Theorem 2.1 of Santos (1996), we have that

$$\mathbb{R}^n = \mathcal{N}(\mathbf{M}\mathbf{A}) \oplus \mathcal{R}(\mathbf{M}\mathbf{A}), \tag{3.5}$$

and

$$\mathbf{x}^k = \mathbf{x}_1^0 + (\mathbf{I} - \mathbf{MA})^k \mathbf{x}_2^0 + [\mathbf{I} - (\mathbf{I} - \mathbf{MA})^k](\mathbf{MA})_2^{-1}(\mathbf{Mb})_2 + k(\mathbf{Mb})_1, \qquad (3.6)$$

where $\mathbf{x}_1^0, (\mathbf{Mb})_1 \in \mathcal{N}(\mathbf{MA})$ and $\mathbf{x}_2^0, (\mathbf{Mb})_2 \in \mathcal{R}(\mathbf{MA})$, are such that

$$\mathbf{Mb} = (\mathbf{Mb})_1 + (\mathbf{Mb})_2, \quad \mathbf{x}^0 = \mathbf{x}_1^0 + \mathbf{x}_2^0, \quad \text{and} \quad (\mathbf{MA})_2 = \mathbf{MA}\Big|_{\mathcal{R}(\mathbf{MA})}.$$

But $(\mathbf{Mb})_1 = 0$, because the method is convergent for every $\mathbf{x}^0$, thus $\mathbf{Mb} \in \mathcal{R}(\mathbf{MA})$, for every $m$-vector $\mathbf{b}$. Therefore, by (1.7) we have that

$$\mathbf{A}(k)(\mathbf{b}) = \mathbf{A}[\mathbf{I} - (\mathbf{I} - \mathbf{MA})^k](\mathbf{MA})_2^{-1}\mathbf{Mb} + \mathbf{Ax}_1^0 + \mathbf{A}(\mathbf{I} - \mathbf{MA})^k \mathbf{x}_2^0. \qquad (3.7)$$

Now

$$\begin{aligned}
\mathbf{A}[\mathbf{I} - (\mathbf{I} - \mathbf{MA})^k](\mathbf{MA})_2^{-1}\mathbf{M} &= -\mathbf{A}\sum_{j=1}^k \binom{k}{j}(-\mathbf{MA})^j(\mathbf{MA})_2^{-1}\mathbf{M} \\
&= -\sum_{j=1}^k \binom{k}{j}(-\mathbf{AM})^j \mathbf{A}(\mathbf{MA})_2^{-1}\mathbf{M} \\
&= -\sum_{j=1}^k \binom{k}{j}(-\mathbf{AM})^{j-1}[-\mathbf{A}(\mathbf{MA}(\mathbf{MA})_2^{-1})\mathbf{M}] \\
&= -\sum_{j=1}^k \binom{k}{j}(-\mathbf{AM})^{j-1}(-\mathbf{AM}) \\
&= \mathbf{I} - (\mathbf{I} - \mathbf{AM})^k. \qquad (3.8)
\end{aligned}$$

Replacing (3.8) in (3.7) the result follows.                                  □

We observe that, in spite of the fact that the method is linear and stationary, if the starting point is not zero, the influence operator is not linear but affine and in many situations (see Herman 1980; Kaufman 1987) it is convenient to choose $x^0 \neq 0$. We will return to this point in the next section.

The expression (3.3) contains all convergent stationary methods that are consistent with the system $\mathbf{MA}x = \mathbf{Mb}$. Important particular cases are the methods considered by Wahba (1987) of the form (3.1) where $\mathbf{M} = \tilde{\mathbf{M}}\mathbf{A}^t$. Section 4 presents numerical experiments with ART, that also fits in Theorem 2 general setting as we will show later.

**Proposition 1.**   *For a method of the form (3.3), the following equality holds,*

$$Tr(\mathbf{I}_m - \mathbf{A}_0(k)) = \frac{1}{n}Tr[(m - n)\mathbf{I}_n + n(\mathbf{I}_n - \mathbf{MA})^k], \qquad (3.9)$$

*where $\mathbf{A}_0(k)$ is the linear part of the influence operator $\mathbf{A}(k)$.*

   *Proof:*   By Theorem 2,

$$\mathbf{A}_0(k) = \mathbf{I}_m - (\mathbf{I}_m - \mathbf{AM})^k = -\sum_{j=1}^k \binom{k}{j}(-\mathbf{AM})^j. \qquad (3.10)$$

Then, using the fact that $Tr(\mathbf{AB}) = Tr(\mathbf{BA})$ and (3.10) we deduce that

$$Tr(\mathbf{A}_0(k)) = Tr[-\sum_{j=1}^{k}\binom{k}{j}(-\mathbf{MA})^j] = Tr[\mathbf{I}_n - (\mathbf{I}_n - \mathbf{MA})^k]. \tag{3.11}$$

Using (3.11), we obtain

$$Tr(\mathbf{I}_m - \mathbf{A}_0(k)) = m - Tr(\mathbf{A}(k)) = m - n + Tr[(\mathbf{I}_n - \mathbf{MA})^k]. \tag{3.12}$$

and (3.9) follows. □

Thus,

$$\frac{1}{m}Tr(\mathbf{I} - \mathbf{A}_0(k)) \approx \frac{\mathbf{w}^t((m-n)\mathbf{w} - n(\mathbf{I}_n - \mathbf{MA})^k\mathbf{w})}{m\mathbf{w}^t\mathbf{w}},$$

where $(\mathbf{I}_n - \mathbf{MA})^k\mathbf{w}$ is equal to the $k$th iteration of (3.3), with $\mathbf{x}^0 = \mathbf{w}$ and $\mathbf{b} = 0$ and $\mathbf{w} = (w_1, \ldots, w_n)^t \in \mathbb{R}^n$ is a random vector from a normal distribution with standard deviation 1, that is, $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{n \times n})$.

From Proposition 1 we deduce the following algorithm to compute for each $k$ the randomized GCV function.

**Algorithm 1.**
   (i) Generate a pseudo-random vector $\mathbf{w} = (w_1, \ldots, w_n)^t \in \mathbb{R}^n$ from a normal distribution with standard deviation 1, that is, $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{n \times n})$;
   (ii) take

$$\Phi(k) = \left(\frac{\mathbf{w}^t((m-n)\mathbf{w} - n(\mathbf{I}_n - \mathbf{MA})^k\mathbf{w})}{m\mathbf{w}^t\mathbf{w}}\right)^2. \tag{3.13}$$

   where $(\mathbf{I}_n - \mathbf{MA})^k\mathbf{w}$ is equal to iteration $k$ of (3.3), with $x^0 = \mathbf{w}$ and $\mathbf{b} = 0$. The expression (3.13) is the approximation of $(\frac{1}{m}Tr[\mathbf{I}_m - \mathbf{A}_0(k)(\mathbf{b})])^2$;
   (iii) finally, compute

$$\frac{\frac{1}{m}||\mathbf{b} - \mathbf{Ax}^k||^2}{\Phi(k)}, \tag{3.14}$$

   that approximates $V(k)$.

In this way we save a product of the matrix $\mathbf{A}$ by a vector, an important saving in very large-scale problems like tomography, where this is the main computational cost.

## 4. A RANDOMIZED ESTIMATE OF THE TRACE OF GENERAL MATRICES

We present some results using a Monte Carlo type method in order to obtain a reliable estimate for the trace. Our results follow closely those given by Girard (1989) for symmetric matrices.

**Theorem 3.** *Let $\mathbf{w} = (w_1, \ldots, w_n)^t$ be a vector of random components with normal distribution, that is, $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_{n \times n})$. Let $\mathbf{B}$ an $n \times n$ matrix, and $T_{\mathbf{B}}(\mathbf{w})$ the random variable*

$$T_{\mathbf{B}}(\mathbf{w}) = \mathbf{w}^t \mathbf{B} \mathbf{w} \,. \tag{4.1}$$

*Then $T_{\mathbf{B}}(\mathbf{w})$ is an unbiased estimator of $Tr(\mathbf{B})$ with standard deviation $(Tr(\mathbf{B}\mathbf{B}^t) + Tr(\mathbf{B}^2))^{1/2}$, that is,*

$$ET_{\mathbf{B}}(\mathbf{w}) = Tr(\mathbf{B}), \tag{4.2}$$

*and*

$$\sigma(T_{\mathbf{B}}(\mathbf{w})) = (Tr(\mathbf{B}\mathbf{B}^t) + Tr(\mathbf{B}^2))^{1/2} \,. \tag{4.3}$$

**Proof:** Taking into account that $E(w_i^2) = 1$ and $E(w_i w_j) = 0$, if $i \neq j$, then

$$ET_B(\mathbf{w}) = E\left(\sum_{i,j} B_{ij} w_i w_j\right) = \sum_{i,j} B_{ij} E(w_i w_j) = Tr(\mathbf{B}) \,. \tag{4.4}$$

Now, from the fact that $E(w_i^4) = 3$, $E(w_i w_j w_k w_l) = 0$, if $i \neq j, k, l$ and $E(w_i^2 w_j^2) = 1$, if $i \neq j$, we deduce that

$$
\begin{aligned}
E((T_{\mathbf{B}}(\mathbf{w}))^2) &= E\left(\left(\sum_{i,j} B_{ij} w_i w_j\right)^2\right) \\
&= 3\sum_i B_{ij}^2 + \sum_{i \neq j} B_{ii} B_{jj} + \sum_{i \neq j} B_{ij} B_{ij} + \sum_{i \neq j} B_{ij} B_{ji} \\
&= Tr(\mathbf{B})^2 + Tr(\mathbf{B}\mathbf{B}^t) + Tr(\mathbf{B}^2) \,.
\end{aligned}
\tag{4.5}
$$

Then,

$$
\begin{aligned}
\sigma^2(T_{\mathbf{B}}(\mathbf{w})) &= E((T_{\mathbf{B}}(\mathbf{w}))^2) - E(T_{\mathbf{B}}(\mathbf{w}))^2 \\
&= Tr(\mathbf{B})^2 + Tr(\mathbf{B}\mathbf{B}^t) + Tr(\mathbf{B}^2) - Tr(\mathbf{B})^2 \,,
\end{aligned}
\tag{4.6}
$$

and the result follows.

$\square$

**Corollary 1.** *Let $\mathbf{w}$ and $\mathbf{B}$ be as in the preceding theorem. Let $T_{\mathbf{B}}^*(\mathbf{w})$ be a random variable defined by*

$$T_{\mathbf{B}}^*(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{B} \mathbf{w}}{\mathbf{w}^t \mathbf{w}} \,. \tag{4.7}$$

*Then $T_{\mathbf{B}}^*(\mathbf{w})$ is an unbiased estimator of $\frac{1}{n} Tr(\mathbf{B})$ with standard deviation*

$$\sqrt{\frac{\dfrac{Tr(\mathbf{B}\mathbf{B}^t) + Tr(\mathbf{B}^2)}{n} - 2\left(\dfrac{Tr(\mathbf{B})}{n}\right)^2}{n+2}} \,,$$

*that is,*

$$ET_{\mathbf{B}}^{*}(\mathbf{w}) = \frac{1}{n}Tr(\mathbf{B}),$$ (4.8)

*and*

$$\sigma(T_{\mathbf{B}}^{*}(\mathbf{w})) = \sqrt{\frac{\frac{Tr(\mathbf{BB}^{t}) + Tr(\mathbf{B}^{2})}{n} - 2\left(\frac{Tr(\mathbf{B})}{n}\right)^{2}}{n+2}}.$$ (4.9)

**Proof:** One way to compute the moments of $T_{\mathbf{B}}^{*}(\mathbf{w})$ is by observing that $(\mathbf{w}^{t}\mathbf{Bw})/(\mathbf{w}^{t}\mathbf{w})$ and $\mathbf{w}^{t}\mathbf{w}$ are independently distributed. Then, it is easily deduced that the moments of $T_{\mathbf{B}}^{*}(\mathbf{w})$ are equal to the moments of $\mathbf{w}^{t}\mathbf{Bw}$ divided by the moments of $\mathbf{w}^{t}\mathbf{w}$, as it was done by Girard (1989). So,

$$E\left(\frac{\mathbf{w}^{t}\mathbf{Bw}}{\mathbf{w}^{t}\mathbf{w}}\right) = \frac{E(\mathbf{w}^{t}\mathbf{Bw})}{E(\mathbf{w}^{t}\mathbf{w})} = \frac{Tr(\mathbf{B})}{n},$$ (4.10)

and

$$E\left(\left(\frac{\mathbf{w}^{t}\mathbf{Bw}}{\mathbf{w}^{t}\mathbf{w}}\right)^{2}\right) = \frac{Tr(\mathbf{B})^{2} + Tr(\mathbf{BB}^{t}) + Tr(\mathbf{B}^{2})}{n^{2} + 2n}.$$ (4.11)

Now,

$$\begin{aligned}
\sigma^{2}(T_{\mathbf{B}}^{*}(\mathbf{w})) &= E((T_{\mathbf{B}}^{*}(\mathbf{w}))^{2}) - E(T_{\mathbf{B}}^{*}(\mathbf{w}))^{2} \\
&= \frac{n(Tr(\mathbf{BB}^{t}) + Tr(\mathbf{B}^{2})) - 2Tr(\mathbf{B})^{2}}{(n^{2} + 2n)n} \\
&= \frac{\frac{Tr(\mathbf{BB}^{t}) + Tr(\mathbf{B}^{2})}{n} - 2\left(\frac{Tr(\mathbf{B})}{n}\right)^{2}}{n+2}.
\end{aligned}$$ (4.12)

And the result follows.

$\square$

We observe that, if $\mathbf{B}$ is symmetric, Corollary 1 corresponds to Theorem 2.2 of Girard (1989).

If $||\mathbf{B}||_{1}$, $||\mathbf{B}||_{2}$, or $||\mathbf{B}||_{\infty}$, remain bounded when $n$ tends to infinite, then $\sigma(T_{B}^{*}(\mathbf{w}))$ tends to zero at a rate $1/\sqrt{n+2}$ as proven in the following. This is a reasonable assumption, because the matrices for which we compute the trace can be viewed as discretization of continuous operators in some Banach space.

**Corollary 2.** *Let $\mathbf{w}$ and $\mathbf{B}$ be as in Theorem 3. If $||\mathbf{B}||_{i} \leq \mathbf{M}$, for $i = 1, 2$ or $\infty$ and some $\mathbf{M} > 0$, then*

$$\sigma(T_{\mathbf{B}}^{*}(\mathbf{w})) \leq \frac{2\mathbf{M}}{\sqrt{n+2}}.$$ (4.13)

**Proof:**    The following inequalities are valid (see, e.g., Golub and Van Loan 1996, chap. 2)

$$\max |\mathbf{A}_{ij}| \leq ||\mathbf{A}||_i, \tag{4.14}$$

and

$$||\mathbf{A}||_F \leq \sqrt{n} ||\mathbf{A}||_i; \tag{4.15}$$

for $i = 1, 2$ e $\infty$, then we obtain that

$$|Tr(\mathbf{B})| \leq \sum_i |B_{ii}| \leq n \max |B_{ii}| \leq n ||\mathbf{B}||_i \leq n \mathbf{M}, \tag{4.16}$$

leading to

$$|Tr(\mathbf{B}^2)| \leq n ||\mathbf{B}^2||_i \leq n ||\mathbf{B}||_i^2 \leq n \mathbf{M}^2; \tag{4.17}$$

and

$$Tr(\mathbf{B}\mathbf{B}^t) = \sum_{ij} B_{ij}^2 = ||\mathbf{B}||_F^2 \leq n ||\mathbf{B}||_i^2 \leq n \mathbf{M}^2, \tag{4.18}$$

for $i = 1, 2$ and $\infty$. The result is an immediate consequence of the above inequalities and (4.12).                                                                                      □

Inequality (4.13) is far from being the best possible one. For less general cases, Girard (1989) proved that $\sigma(T_B^*(\mathbf{w}))$ can decrease at a $1/n$ rate.


# 5. AN APPLICATION OF GCV TO POSITRON EMISSION TOMOGRAPHY

The goal of emission computed tomography (ECT) is the quantitative determination of the moment-to-moment changes in the chemistry and flow physiology of radioactive labeled components inside the body. The mathematical problem consists of reconstructing a function representing the distribution of radioactivity in a body cross-section from measured data that are the total activity along lines of known location. One of the main differences between this problem and the one arising in X-ray tomography (Herman 1980) is that here measurements tend to be much more noisy, so direct inversion using convolution backprojection (CBP) does not necessarily give the best results (see Vardi, Shepp, and Kaufman 1985).

In positron emission tomography (PET) (Ter-Pogossian et al. 1980) the isotope used emits positrons which annihilate with nearby electrons generating two photons traveling away from each other in (nearly) opposite directions; the number of such photons pairs (detected in time coincidence) for each line or pair of detectors is related to the integral of the concentration of the isotope along the line.

Suppose now that we discretize the problem by subdividing the reconstruction region into $n$ small square-shaped picture elements (pixels, for short) and we assume that the

activity in each pixel $j$ is a constant, denoted by $x_j$. If we count $b_i$ coincidences along $m$ lines and $a_{ij}$ denotes the probability that a photon emitted by pixel $j$ is detected by pair $i$, then $b_i$ is a sample from a Poisson distribution whose expected value is $\sum_{j=1}^{n} a_{ij}x_j$. Taking this into account, it was suggested by Rockmore and Macovski (1976) to estimate $x = (x_1, \ldots, x_n)^t$ by maximizing the Poisson likelihood. Shepp and Vardi (1982) suggested the use of the EM algorithm (Dempster, Laird, and Rubin 1977) for this maximization, that has since then become very popular in the field. Herman and Meyer (1993) proposed the use of ART (from algebraic reconstruction technique) for PET obtaining a speed increase of at least one order of magnitude with respect to the EM, and comparable image quality. ART (see Herman 1980), that is, the algorithm from Kaczmarz (1937), is defined by

$$\begin{cases} \mathbf{x}^{(k+1,1)} = \mathbf{x}^k \\ \mathbf{x}^{(k+1,i+1)} = \mathbf{x}^{(k+1,i)} + \omega \dfrac{(b_i - \mathbf{a}_i^t x^{(k+1,i)})}{||a_i||^2} a_i \\ \mathbf{x}^{k+1} = \mathbf{x}^{(k+1,m+1)}, \end{cases} \tag{5.1}$$

for $i = 1, \ldots, m$; where $\omega$ is a relaxation parameter in the interval $(0, 2)$ and $\mathbf{A} = [a_1, \ldots, a_m]^t = \{a_{ij}\}$.

It was proven by Björck and Elfving (1979) that the above method corresponds to applying SOR to the system

$$\begin{cases} \mathbf{A}\mathbf{A}^t\mathbf{y} = \mathbf{b} \\ \mathbf{x} = \mathbf{A}^t\mathbf{y} \end{cases}, \tag{5.2}$$

that is,

$$\begin{cases} \mathbf{y}^{k+1} = \mathbf{y}^k + \hat{\mathbf{M}}(\mathbf{b} - \mathbf{A}\mathbf{A}^t\mathbf{y}^k) \\ \mathbf{x}^k = \mathbf{A}^t\mathbf{y}^k \end{cases}. \tag{5.3}$$

Therefore, the method (5.1) can be rewritten as

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \mathbf{A}^t\hat{\mathbf{M}}(\mathbf{b} - \mathbf{A}\mathbf{x}^k).$$

This shows that the ART method fits in the framework of Theorem 2.

The following remark ensures the applicability of Corollary 2 to ART for the PET problem. We make the realistic assumption that $m = pn$ (usually $p = 3$ in 2-D PET, for example).

**Remark.** ART is a sequence of relaxed orthogonal projections onto hyperplanes, so $\mathbf{I} - \mathbf{MA}$ in (3.3) is a product of matrices that are orthogonal projections onto subspaces. Therefore $||\mathbf{I} - \mathbf{MA}||_2 \leq 1$ and by Corollary 2 in the next section, we can apply the Monte-Carlo approach (Algorithm 1) to this problem.

In our numerical experiments we used the programming system SNARK93, developed by the Medical Image Processing Group of the University of Pennsylvania (Browne, Herman, and Odhner 1993). The images to be reconstructed (phantom) were obtained from
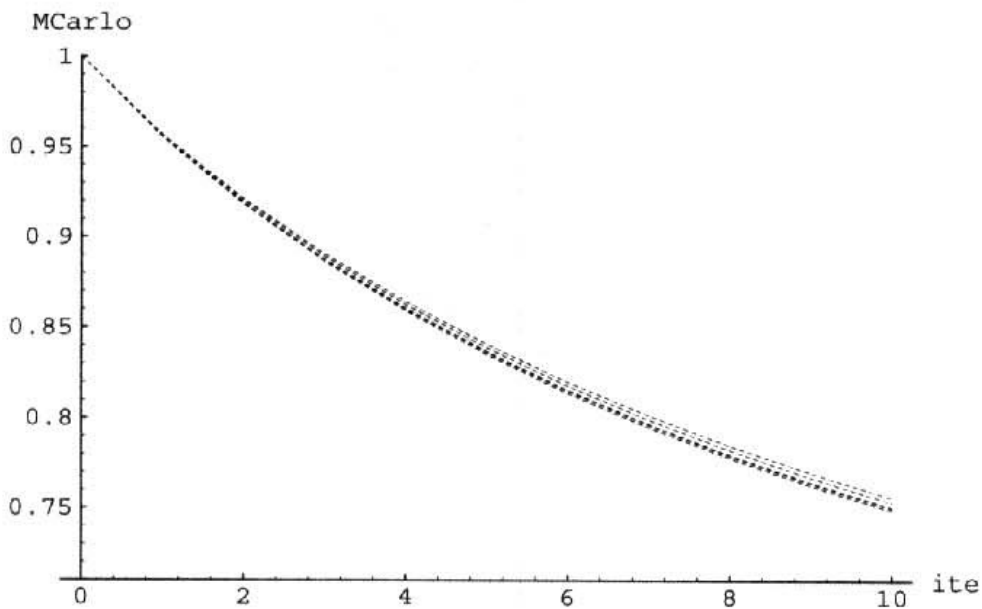
MCarlo



Figure 1.    Five Monte-Carlo estimates of $[\frac{1}{30292}Tr(\mathbf{I}_{30292} - \mathbf{A}(k))]^2$ for ART with $\omega = 0.025$.
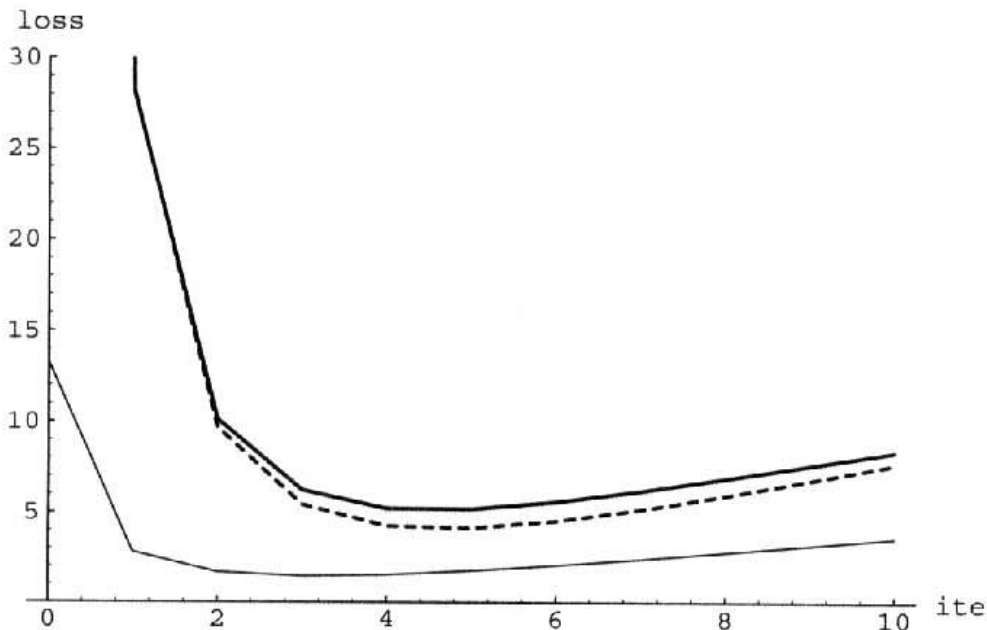
loss



Figure 2.   ART with $\omega = 0.025$, white noise. The thin solid line corresponds to $100\,\frac{1}{95^2}||\mathbf{x}^k - \mathbf{x}||^2$; the thick one to $\frac{1}{30292}||\mathbf{Ax}^k - \mathbf{Ax}||^2$; the dashed line to MCarlo GCV $-100$.
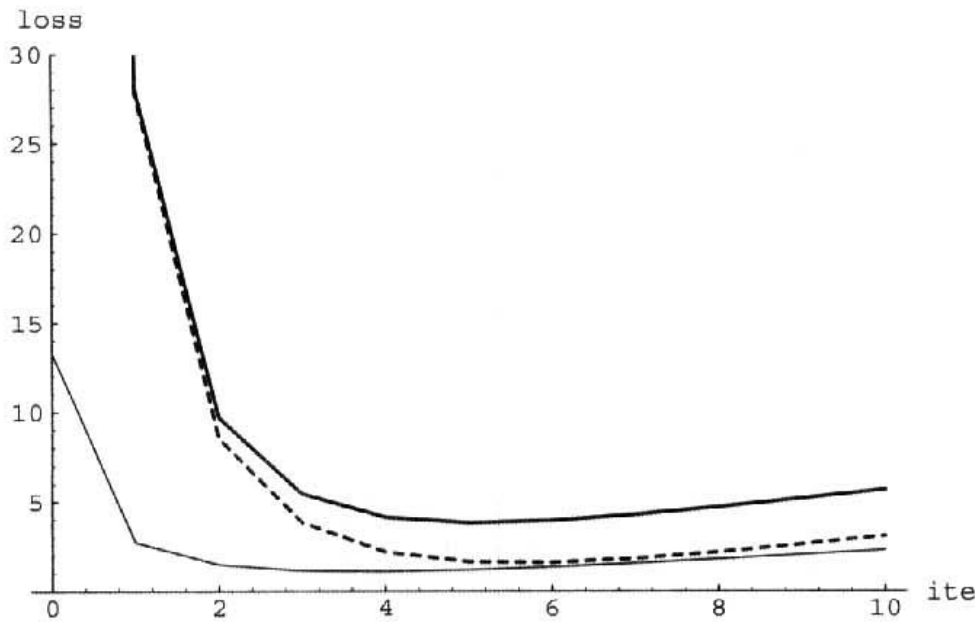
*Figure 3. Loss function for ART with $\omega = 0.025$, PET errors. The thin solid line corresponds to $100\frac{1}{95^2}||\mathbf{x}^k - \mathbf{x}||^2$; the thick one to $\frac{1}{30292}||\mathbf{Ax}^k - \mathbf{Ax}||^2$; the dashed line to MCarlo GCV $-70$.*

a computerized atlas based on typical values inside the brain, as in Herman and Meyer (1993). The data collection geometry was a divergent one, simulating a typical PET data acquisition (Browne, Herman, and Odhner 1993). The equations are taken counter-clockwise inside each view (set of rays associated with one source) and views were also considered counter-clockwise. The order in which rays are used can substantially improve the convergence rate of the algorithm (see Herman and Meyer 1993; van Dijke 1992) but here we are interested only in validating GCV as a stopping rule.

We used a discretization with $n = 95 \times 95$ pixels and the divergent geometry had 300 views, of 101 rays each, a total number of $m = 30{,}292$ equations (8 rays were not considered because their lack of intersection with the image region). The starting point was a uniform image $x^0 = (a, \ldots, a)^t$, where $a$ is an approximation of the average density of the phantom given by

$$a = \frac{\sum_{i=1}^{m} b_i}{\sum_{i=1}^{m}\sum_{j=1}^{n} a_{ij}} \,. \tag{5.4}$$

A uniform nonzero starting point was advocated by Kaufman (1987) and it is widely accepted as the best choice for many reseachers in ECT. For most of the experiments $\omega$ was fixed to 0.025, because we know that for small values of $\omega$, the method approximates a weighted least squares solution (Herman and Meyer 1993).

To show that the Monte Carlo approach provides very good estimates for the denominator of (2.5) we performed several experiments applying the first two steps, (i) and (ii), of Algorithm 1. In those experiments we used several choices for the relaxation parameter $\omega$ and for all of them the results were very similar. Figure 1 shows the results for $\omega = 0.025$.
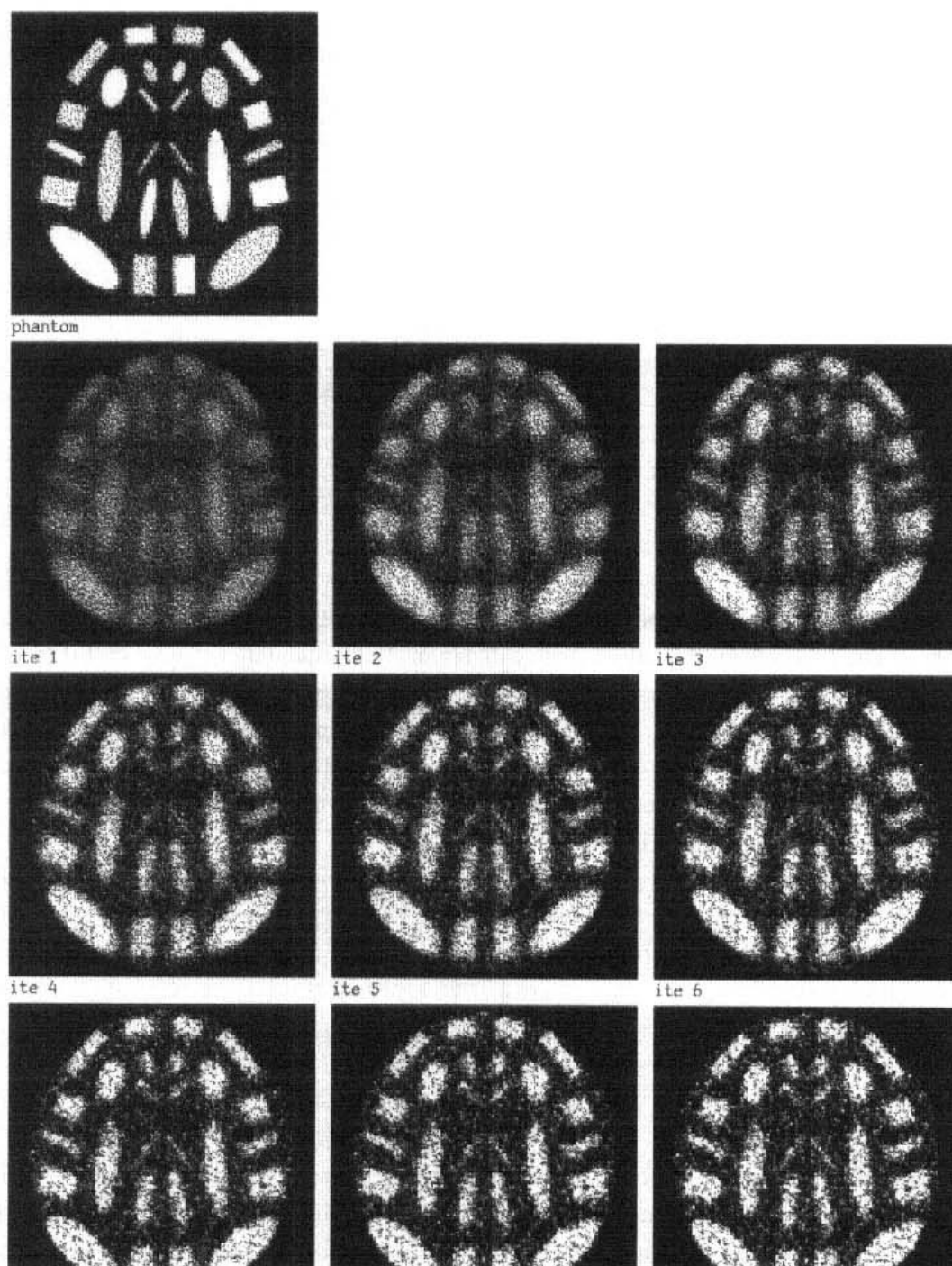
phantom

ite 1            ite 2            ite 3

ite 4            ite 5            ite 6

Figure 4.    A phantom and 9 image iterates corresponding to the plots in Figure 3.

We picked up five random vectors $\mathbf{w}$, and, for each one, we plotted $\Phi(k)$, the estimate of the GCV functional denominator, against the iteration number $k$. As expected, the curves are decreasing and almost coincident.

We performed several experiments applying Algorithm 1 for several choices of the relaxation parameter $\omega$ and for all of them the results were very similar. In Figure 2 we plot the functions $100\frac{1}{95^2}||\mathbf{x}^k - \mathbf{x}||^2$, $\frac{1}{30292}||\mathbf{Ax}^k - \mathbf{Ax}||^2$ and Monte Carlo GCV $- 100$ against the iteration number $k$, for $\omega = 0.025$. The error vectors $\epsilon$ were white noise with standard deviation equal to 10, that is, each component of $\epsilon$ is a pseudo-random number with a normal
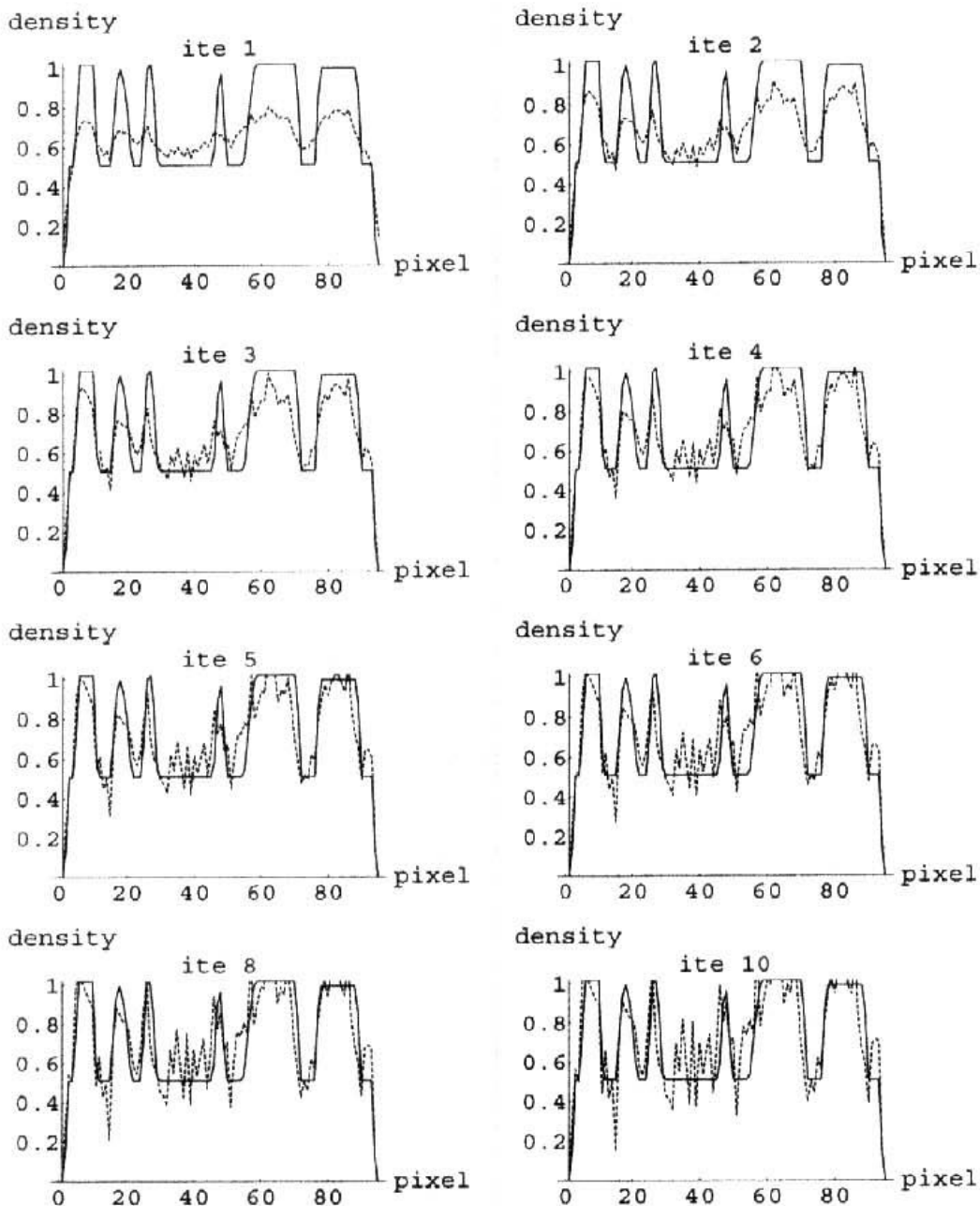


Figure 5. Comparison of exact emission densities for column 42 corresponding to the images in Figure 4.

distribution, zero mean, and standard deviation 10, that is, $\epsilon \sim \mathcal{N}(0, 100\mathbf{I}_{30292})$. As expected the minimum of curve of the Monte Carlo GCV coincides with that of $\frac{1}{30292}||\mathbf{A}\mathbf{x}^k - \mathbf{A}\mathbf{x}||^2$. In Figure 3 we plot the same functions for ART, but in this case we use typical PET errors as in Herman and Odhner (1991), derived from a Poisson distribution; the total photon count number was 2,022,085. The graphs show that, when the errors satisfy the theoretical hypotheses (Figure 2), and also when they are typical PET errors, not only the minimum of Monte Carlo GCV coincides with that of $||\mathbf{A}\mathbf{x}^k - \mathbf{A}\mathbf{x}||^2$, but also the curves are very similar.

Figure 4 shows the images corresponding to several iterations of ART with data as in Figure 3. Figure 5 shows a cross section of line 42 for the images of Figure 4.

## 6. CONCLUDING REMARKS

We have shown in this aricle a method to apply GCV as a stopping rule to general iterative linear stationary methods for large scale ill-posed problems. Our choice for ART as the method and PET as the problem to illustrate this applicability is a consequence of our particular interest in tomography. Nevertheless, it is clear that this article gives the foundations and shows how to apply GCV to LSIM's like SOR, JOR, and so on, when solving large scale linear systems of equations originated in ill-posed problems.

## ACKNOWLEDGMENTS

## REFERENCES

Allen, D. M. (1974), "The Relationship Between Variable Selection and Data Augmentation and a Method for Prediction," *Technometrics*, 16, 125–127.

Björck, A., and Elfving, T. (1979), "Accelerated Projection Methods for Computing Pseudoinverse Solutions of Systems of Linear Equations," *BIT*, 19, 145–163.

Browne, J. A., Herman, G. T., and Odhner, D. (1993), "SNARK93 A Programming System for Image Reconstruction from Projections," Technical Report MIPG198, Department of Radiology, University of Pennsylvania.

Craven, P., and Wahba, G. (1979), "Smoothing Noisy Data With Spline Functions," *Numerische Mathematik*, 31, 377–403.

Dempster, A. P., Laird, N. M., and Rubin, D. D. (1977), "Maximum Likelihood from Incomplete Data via the EM Algorithm," *Journal of the Royal Statistical Society*, Series B, 39, 1–38.

Fleming, H. E. (1990), "Equivalence of Regularization and Truncated Iteration in the Solution of Ill-Posed Image Reconstruction Problems," *Linear Algebra and its Applications*, 130, 133–150.

Girard, D. A. (1989), "A Fast 'Monte-Carlo Cross-Validation' Procedure for Large Least Squares Problems with Noisy Data," *Numerische Mathematik*, 56, 1–23.

Golub, G., and Matt, U. von (1997), "Generalized Cross-Validation for Large-Scale Problems," *Journal of Computational and Graphical Statistics*, 6, 1–34.

Golub, G. H., Heath, M. T., and Wahba, G. (1979), "Generalized Cross-Validation as a Method for Choosing a Good Ridge Parameter," *Technometrics*, 21, 215–223.

Golub, G. H., and Van Loan, C. F. (1996), *Matrix Computations*, Baltimore, MD: Johns Hopkins.

Herman, G. T. (1980), *Image Reconstruction from Projections: The Fundamentals of Computerized Tomography*, New York: Academic Press.

Herman, G. T., and Meyer, L. B. (1993), "Algebraic Reconstruction Techniques can be Made Computationally Efficient," *IEEE Transactions on Medical Imaging*, 12, 600–609.

Herman, G. T., and Odhner, D. (1991), "Performance Evaluation of an Iterative Image Reconstruction Algorithm for Positron Emission Tomography," *IEEE Transactions on Medical Imaging*, 10, 336–346.

Kaczmarz, S. (1937), "Angenährte Auflösung von Systemen linearer Gleichungen," *Bulletin de l'Académie Polonaise des Sciences et Lettres A*, 35, 355–357.

Kaufman, L. (1987), "Implementing and Accelerating the EM Algorithm for Positron Emission Tomography," *IEEE Transactions on Medical Imaging*, 6, 37–51.

Perry, K. M., and Reeves, S. J. (1994), "A Practical Stopping Rule for Iterative Signal Restoration," *IEEE Transactions on Signal Processing*, 42, 1829–1833.

Reeves, S. J. (1992), "A Cross-Validation Framework for Solving Imaging Restoration Problems," *Journal of Visual Communication and Image Representation*, 3, 433–445.

——— (1994), "Optimal Space-Varying Regularization in Iterative Image Restoration," *IEEE Transactions on Image Processing*, 3, 319–324.

Rockmore, A., and Macovski, A. (1976), "A Maximum Likelihood Approach to Emission Image Reconstruction from Projections," *IEEE Transactions on Nuclear Science*, 23, 1428–1432.

Santos, R. J. (1996), "Equivalence of Regularization and Truncated Iteration for General Ill-Posed Problems," *Linear Algebra and its Applications*, 236, 25–33.

Shepp, L., and Vardi, Y. (1982), "Maximum Likelihood Reconstruction for Emission Tomography," *IEEE Transactions on Medical Imaging*, MI-1, 113–121.

Silvey, S. D. (1970), *Statistical Inference*, Harmondsworth: Penguin.

Ter-Pogossian, M. M., et al. (1980), "Positron Emission Tomography," *Scientific American*, October, pp. 170–181.

van der Sluis, A., and van der Vorst, H. A. (1987), "Numerical Solution of Large, Sparse Linear Algebraic Systems Arising from Tomographic Problems," in *Seismic Tomography*, ed. G. Nolet, Dordrecht: Reidel, 49–83.

——— (1990), "SIRT and CG Type Methods for the Iterative Solution of Sparse Linear Least Squares Problems," *Linear Algebra and its Applications*, 130, 257–303.

van Dijke, M. (1992), *Iterative Methods in Image Reconstruction*, PhD thesis, Rijksuniversiteit Utrecht, Utrecht, The Netherlands.

Vardi, Y., Shepp, L. A., and Kaufman, L. (1985), "A Statistical Model for Positron Emission Tomography," *Journal of the American Statistical Association*, 80, 8–37.

Wahba, G. (1987), "Three Topics in Ill-Posed Problems," in *Inverse and Ill-Posed Problems*, eds. H. Engl and C. Groetsch, New York: Academic Press, 37–51.