# Exploring Estimator Bias-Variance Tradeoffs Using the Uniform CR Bound

Alfred O. Hero, III, *Member, IEEE*, Jeffrey A. Fessler, *Member, IEEE*, and Mohammad Usman, *Member, IEEE*

*Abstract*—We introduce a plane, which we call the delta–sigma plane, that is indexed by the norm of the estimator bias gradient and the variance of the estimator. The norm of the bias gradient is related to the maximum variation in the estimator bias function over a neighborhood of parameter space. Using a uniform Cramér–Rao (CR) bound on estimator variance, a delta–sigma tradeoff curve is specified that defines an "unachievable region" of the delta–sigma plane for a specified statistical model. In order to place an estimator on this plane for comparison with the delta–sigma tradeoff curve, the estimator variance, bias gradient, and bias gradient norm must be evaluated. We present a simple and accurate method for experimentally determining the bias gradient norm based on applying a bootstrap estimator to a sample mean constructed from the gradient of the log-likelihood. We demonstrate the methods developed in this paper for linear Gaussian and nonlinear Poisson inverse problems.

## I. INTRODUCTION

**T**HE goal of this work is to quantify fundamental tradeoffs between the bias and variance functions for parametric estimation problems. Let $\underline{\theta} = [\theta_1, \cdots, \theta_n]^T \in \Theta$ be a vector of unknown and nonrandom parameters that parameterize the density $f_Y(y; \underline{\theta})$ of an observed random variable $Y$. The parameter space $\Theta$ is assumed to be an open subset of $n$-dimensional Euclidean space $\mathbb{R}^n$. For fixed $\underline{\theta}$, let $\hat{t} = \hat{t}(Y)$ be an estimator of the scalar $t_\theta$, where $t: \Theta \to \mathbb{R}$ is a specified function. Let this estimator have bias $b_\theta = E_\theta[\hat{t}] - t_\theta$ and variance $\sigma_\theta^2 = E_\theta[(\hat{t} - t_\theta)^2]$. Bias is due to 'mismatch' between the average value of the estimator and the true parameter, whereas variance arises from fluctuations in the estimator due to statistical sampling.

In most applications, estimator designs are subject to a tradeoff between bias and variance. For example, in nonparametric spectrum estimation [1], smoothing methods have long been used to reduce the variance of the periodogram at the expense of increased bias [2], [3]. In image restoration, regularization is frequently implemented to reduce noise amplification (variance) at the expense of reduced spatial resolution (bias) [4]. In multiple regression with multicollinearity, biased shrinkage

estimators [5] and biased ridge estimators [6] are used to reduce variance of the ordinary least squares estimator. The quantitative study of estimator bias and variance has been useful for characterizing statistical performance for many statistical signal processing applications including tomographic reconstruction [7]–[9], functional imaging [10], nonlinear and morphological filtering [11], [12], and spectral estimation of time series [13], [14].

However, the plane parameterized by the bias and variance $b_\theta$ and $\sigma_\theta^2$ is not useful for studying fundamental tradeoffs since an estimator can always be found that makes both the bias and variance zero at a given point $\underline{\theta}$. Furthermore, use of bias can be misleading: Even a very large bias is removable if it is constant. In this work, we consider the plane parameterized by the norm or length of the bias gradient $\delta_{\underline{\theta}} = \|\nabla b_{\underline{\theta}}\|$ and the square root variance $\sqrt{\sigma_{\underline{\theta}}^2}$, which we call the delta-sigma or $\delta\sigma$ plane. The norm of the bias gradient is directly related to the maximal variation of the bias function over a neighborhood of $\underline{\theta}$ induced by the norm and is unaffected by constant estimator bias components. By appropriate choice of norm, the bias gradient length can be related to the overall bias variation over any prior ellipsoidal region of parameter values. For the inverse problems studied here, we select the norm to correspond to an *a priori* smoothness constraint on the object.

This paper provides a means for specifying unachievable regions in the $\delta\sigma$ plane via fundamental delta–sigma tradeoff curves. These curves are generated using an extension of the Cramér-Rao (CR) lower bound on the variance of biased estimators presented in [15]. This extension is called the uniform CR bound. In [15], the bound was derived only for an unweighted Euclidean norm on the bias gradient and for nonsingular Fisher information. Therein, the reader was cautioned that the resulting bound will generally depend on the units and dimensions used to express each of the parameters. It was also pointed out in [15] that the user should identify an ellipsoid of expected parameter variations, which will depend in the user's units, and perform a normalizing transformation of the ellipsoid to a spheroid prior to applying the bound. This parameter transformation is equivalent to using a diagonally weighted bias gradient norm constraint in the original untransformed parameter space. The uniform CR bound presented in this paper generalizes [15] to allow functional estimation, to cover the case of singular or ill-conditioned Fisher matrices, and to account for a general norm constraint on the bias gradient. Some elements of the latter generalization were first presented in [16].

The methods described herein can be used for system optimization, i.e., to choose the system that minimizes the size of the unachievable region when estimator unbiasedness is an overly stringent or unrealistic constraint [17], or they can be used to gauge the closeness to optimality of biased estimators in terms of their nearness to the unachievable region [18]. Alternatively, as discussed in more detail in [15], these results can be used to investigate the reliability of unbiased CR bound studies when small estimator biases may be present. Finally, these results can be used for validation of estimator simulations by empirically verifying that the simulations do not place estimator performance in the unachievable region of the $\delta\sigma$ plane.

In order to place an estimator on the $\delta\sigma$ plane, we must calculate estimator variance and bias gradient norm. For most nonlinear estimators, analytical computation of these quantities is intractable. We present a methodology for experimentally determining these quantities that use the gradient of the log-likelihood function $\nabla \ln f_Y(y;\underline{\theta})$ and a bootstrap-type estimator to estimate the bias gradient norm.

We illustrate these methods for linear Gaussian and nonlinear Poisson inverse problems. Such problems arise in image restoration, image reconstruction, and seismic deconvolution, to name but a few examples. Note that even for the linear Gaussian problem, there may not exist unbiased estimators when the system matrix is ill conditioned or rank deficient [19]. For each model, we compare the performance of quadratically penalized maximum likelihood estimators to the fundamental delta-sigma tradeoff curve. We show that the bias gradient $\nabla b_{\underline{\theta}}$ of these estimators is closely related to the point spread function of the estimator when one wishes to estimate a single component $t_{\underline{\theta}} = \theta_k$. For the full-rank linear Gaussian case, the quadratically penalized likelihood estimator achieves the fundamental delta-sigma tradeoff in the $\delta\sigma$ plane when the roughness penalty matrix is matched to the norm chosen by the user to measure bias gradient length. In this case, the bias gradient norm constraint is equivalent to a constraint on bias variation over a roughness constrained neighborhood of $\underline{\theta}$. We thus have a very strong optimality property: The penalized maximum likelihood estimator minimizes variance over all estimators whose maximal bias variation is bounded over the neighborhood. For the rank-deficient linear Gaussian problem, the uniform CR bound is shown to be achievable by a different estimator under certain conditions. Finally, for the nonlinear Poisson case, an asymptotic analysis shows that the penalized maximum likelihood estimator of [20] achieves the fundamental delta-sigma tradeoff curve for sufficiently large values of the regularization parameter and a suitably chosen penalty matrix. We present simulation results that empirically validate our asymptotic analysis.

## A. Variance, Bias, and Bias Gradient

Let $\hat{t}$ be an estimator of the scalar differentiable function $t_{\underline{\theta}}$. The mean-square error (MSE) is a widely used measure of performance for an estimator $\hat{t}$ and is simply related to the estimator bias $b_{\underline{\theta}}$ and the estimator variance $\sigma_{\underline{\theta}}^2$ through the relation $\mathrm{MSE}_{\underline{\theta}} = \sigma_{\underline{\theta}}^2 + b_{\underline{\theta}}^2$. While the MSE criterion is of

value in many applications, the estimator bias and estimator variance provide a more complete picture of performance than the MSE alone. From $b_{\underline{\theta}}$ and $\sigma_{\underline{\theta}}^2$, one can derive other important measures such as signal-to-noise-ratio $\mathrm{SNR} = |t_{\underline{\theta}} + b_{\underline{\theta}}|^2/\sigma_{\underline{\theta}}^2$, coefficient of variation 1/SNR, and generalized MSE $= \alpha g_1(b_{\underline{\theta}}) + (1 - \alpha)g_2(\sigma_{\underline{\theta}})$, where $\alpha \in [0,1]$ and $g_1, g_2$ are nonnegative functions. The generalized MSE has been used in response surface design [21] and in minimum bias and variance estimation for nonlinear regression models [22], [23]. Furthermore, since they jointly specify the first two moments of the estimator probability distribution, the pair $(b_{\underline{\theta}}, \sigma_{\underline{\theta}}^2)$ provides essential information for constructing and evaluating $\hat{t}$-based hypothesis tests and confidence intervals. Indeed, the popular jacknife method was originally introduced in [24] and [25] to estimate bias and variance of a statistic and to test whether the statistic has prespecified mean [26].

An estimator $\hat{t}$ whose bias function $b: \Theta \to \mathbb{R}$ is constant is as good as unbiased since the bias can be removed without knowledge of $\underline{\theta}$. Therefore, when one is interested in fundamental tradeoffs, it is the bias variation that will be of interest. When the density function $f_Y(y;\underline{\theta})$ is sufficiently smooth to guarantee existence of the Fisher information matrix (which is defined below), $b_{\underline{\theta}}$ is always differentiable, regardless of the form of the estimator, as long as $E_{\underline{\theta}}[\hat{t}^2]$ is upper bounded [27, Lemma 7.2]. In this case, the bias gradient $\nabla b: \Theta \to \mathbb{R}^n$ uniquely specifies the bias $b_{\underline{\theta}}$ up to an additive constant

$$b_{\underline{\theta}} = \sum_{k=1}^{n} \int_{\theta_k^o}^{\theta_k} \frac{\partial b_{\underline{u}}}{\partial u_k} du_k + b_{\underline{\theta}^o}$$

where $\underline{\theta}^o$ is a point such that the line segment connecting $\underline{\theta}^o$ and $\underline{\theta}$ is contained in $\Theta$—such a point is guaranteed to exist when $\Theta$ is convex or star shaped about a point $\underline{\theta}$. Thus the gradient function $\nabla b_{\underline{\theta}} = [\partial b_{\underline{\theta}}/\partial\theta_1, \cdots, \partial b_{\underline{\theta}}/\partial\theta_1]^T$ (which is a *column* vector) characterizes the unremovable bias component of the bias function.

*1) Bias Gradient Norm and Maximal Bias:* Define the norm or length of the bias gradient vector

$$\delta_{\underline{\theta}} = ||\nabla b_{\underline{\theta}}||_C \tag{1}$$

where the norm $|| \cdot ||_C$ is is defined in terms of a symmetric positive definite matrix $C$

$$||\underline{u}||_C^2 = \underline{u}^T C \underline{u}. \tag{2}$$

We will use the notation $||\underline{u}||_2$ to denote the standard Euclidean norm obtained when $C = I$.

The norm of the bias gradient at a point $\underline{u} = \underline{\theta}$ is a measure of the sensitivity of the estimator mean $m_{\underline{u}} = E_{\underline{u}}[\hat{t}]$ to changes in $\underline{u}$ over a neighborhood of $\underline{\theta}$. Below, we derive a relation between bias gradient norm and maximal bias variation over an ellipsoidal neighborhood.

Define the ellipsoidal region of parameter variations $\mathcal{C} = \mathcal{C}(\underline{\theta}, C) = \{\underline{u}: (\underline{u} - \underline{\theta})^T C^{-1}(\underline{u} - \underline{\theta}) \le 1\}$, where $\underline{\theta}$ is a point in $\Theta$, and $C$ is a symmetric positive definite matrix. The maximal width of the ellipsoid is $2\sqrt{\lambda_M^C}$, where $\lambda_M^C$ is the maximal eigenvalue of $C$. Assume that the bias function $b_{\underline{u}}$ is continuously twice differentiable and that the magnitude of

the eigenvalues of the Hessian matrix $\nabla^2 b_{\underline{u}} = \nabla\nabla^T b_{\underline{u}}$ are upper bounded over $\underline{u} \in C$ by a nonnegative constant $\alpha < \infty$. Then, using (1) and the Taylor expansion with remainder, the maximal squared variation of the bias $b_{\underline{u}}$ over $C$ is

$$\max_{\underline{u}\in C} |b_{\underline{u}} - b_{\underline{\theta}}|^2 = \max_{\underline{u}\in C} |\nabla^T b_{\underline{\theta}}\Delta\underline{u} + \tfrac{1}{2}\Delta\underline{u}^T \nabla^2 b_{\underline{\xi}}\Delta\underline{u}|^2 \quad (3)$$

where $\Delta\underline{u} = \underline{u} - \underline{\theta}$, and $\underline{\xi}$ is a point along the line segment joining $\underline{\theta}$ and $\underline{u}$. Now, expanding the square on the right-hand side of (3) and collecting terms, we obtain

$$\max_{\underline{u}\in C} |b_{\underline{u}} - b_{\underline{\theta}}|^2 = \max_{\underline{u}\in C} |\nabla^T b_{\underline{\theta}}\Delta\underline{u}|^2(1 + \epsilon) \quad (4)$$

where $|\epsilon| \leq \rho(1 + 0.25\rho)$ and $\rho = \lambda_M^C \alpha / \sqrt{\nabla^T b_{\underline{\theta}} C \nabla b_{\underline{\theta}}}$. Defining $\Delta\tilde{\underline{u}} = C^{-(1/2)}\Delta\underline{u}$ and using the Cauchy–Schwarz inequality, we obtain

$$\max_{\underline{u}\in C} |\nabla^T b_{\underline{\theta}}\Delta\underline{u}|^2 = \max_{||\Delta\tilde{\underline{u}}||\leq 1} |\nabla^T b_{\underline{\theta}} C^{1/2}\Delta\tilde{\underline{u}}|^2$$

$$= \max_{\Delta\tilde{\underline{u}}} \frac{|\nabla^T b_{\underline{\theta}} C^{1/2}\Delta\tilde{\underline{u}}|^2}{||\Delta\tilde{\underline{u}}||_2}$$

$$= \nabla^T b_{\underline{\theta}} C \nabla b_{\underline{\theta}}. \quad (5)$$

Therefore, combining (5)–(3)

$$\max_{\underline{u}\in C} |b_{\underline{u}} - b_{\underline{\theta}}|^2 = ||\nabla b_{\underline{\theta}}||_C^2(1 + \epsilon). \quad (6)$$

Hence, we see that when $\rho \ll 1, 1 + \epsilon \approx 1$, and the norm $||\nabla b_{\underline{u}}||_C$ is approximately equal to the maximal bias variation over the ellipsoidal neighborhood $C(\underline{\theta}C)$ of $\underline{\theta}$. Note that this occurs when the product of the ellipsoid width $\lambda_M^C$ and the ratio of the curvature $\alpha$ of the bias function to the bias gradient norm $\sqrt{\nabla^T b_{\underline{\theta}} C \nabla b_{\underline{\theta}}}$ is small. For the special case where the bias is a linear function $(b_{\underline{\theta}} = \underline{L}^T\underline{\theta} - c)$ then $\rho = 0$, in which case the relation (6) between bias gradient norm (1) and maximal bias variation (3) is exact.

The above discussion suggests that the choice of norm $||\cdot||_C$ should reflect the range $C$ of joint parameter variations that are of interest to the user. This will be illustrated in Section IV.

## II. UNACHIEVABLE REGIONS

For any estimator with bias gradient norm $\delta_{\underline{\theta}}$ and variance $\sigma_{\underline{\theta}}^2$, we plot the pair $(\delta_{\underline{\theta}}, \sigma_{\underline{\theta}})$ as a coordinate in the plane $\mathbb{R}^2$. We will call this parameterization of the plane the *delta-sigma* or $\delta\sigma$ plane. A region of the $\delta\sigma$ plane is called unachievable if no estimator can exist having coordinates in this region. While no nonempty unachievable region can exist in the bias-variance plane parameterized by $(b_{\underline{\theta}}, \sigma_{\underline{\theta}})$, we will show that interesting unachievable regions almost always exist in the delta-sigma plane.

### A. The Biased CR Bound

The CR lower bound on estimator variance, which was first published by Frechet [28] and later by Darmois [29], Cramer [30], and Rao [31], is commonly used to bound the variance of unbiased estimators. For a biased estimator $\hat{t}$ of $t_{\underline{\theta}}$ with mean

$m_{\underline{\theta}} = E_{\underline{\theta}}[\hat{t}]$, the CR bound has the following form, which is called the *biased CR bound*:

$$\sigma_{\underline{\theta}}^2 \geq [\nabla m_{\underline{\theta}}]^T F_Y^+[\nabla m_{\underline{\theta}}]$$

$$= [\nabla t_{\underline{\theta}} + \nabla b_{\underline{\theta}}]^T F_Y^+[\nabla t_{\underline{\theta}} + \nabla b_{\underline{\theta}}] \quad (7)$$

where $F_Y = F_Y(\underline{\theta})$ is the $n \times n$ Fisher information matrix

$$F_Y = E_{\underline{\theta}}\{[\nabla_{\underline{\theta}} \ln f_Y(Y;\underline{\theta})][\nabla_{\underline{\theta}} \ln f_Y(Y;\underline{\theta})]^T\},$$

and $F_Y^+$ denotes the Moore–Penrose pseudo-inverse matrix of the possibly singular matrix $F_Y$.

The nonsingular-$F_Y$ form of the biased CR bound has been around for some time, e.g., [32]. The more general pseudo-inverse-$F_Y$ form given in (7) is less well known but can be easily derived by identifying $U = \hat{t} - t_{\underline{\theta}}$ and $V = \nabla_{\underline{\theta}} \ln f_Y(Y;\underline{\theta})$ in the relation [33, Lemma 1]

$$E_{\underline{\theta}}[UU^T] \geq E_{\underline{\theta}}[UV^T](E_{\underline{\theta}}[VV^T])^+ E_{\underline{\theta}}[VU^T],$$

and using the identities $E_{\underline{\theta}}[\nabla_{\underline{\theta}} \ln f_Y(Y;\underline{\theta})] = \underline{0}$ and $E_{\underline{\theta}}[\hat{t}\nabla_{\underline{\theta}} \ln f_Y(Y;\underline{\theta})] = \nabla m_{\underline{\theta}}$ (which are easily derivable from (19) below).

The biased CR bound (7) only applies to the class of estimators $\hat{t}$ that have a particular bias gradient function $\nabla b_{\underline{\theta}}$. Therefore, (7) cannot be used to simultaneously bound the variance of several estimators, each of which have different but comparable bias gradients.

### A. The Uniform CR Bound

In [15], a "uniform" CR bound was presented as a way to study the reliability of the unbiased CR bound under conditions of very small estimator bias. In [34], this uniform bound was used to trace out curves over the sigma–delta plane, which includes both large and small biases. The following theorem extends the results of [15] and [34] to allow singular Fisher information matrices, arbitrary weighted Euclidean norm $||\cdot||_C$, and arbitrary differentiable function $t_{\underline{\theta}}$. For a proof of this theorem, see Appendix A.

*Theorem 1:* Let $\hat{t}$ be an estimator of the scalar differentiable function $t_{\underline{\theta}}$ of the parameter $\underline{\theta} = [\theta_1, \cdots, \theta_n]^T$. For a fixed $\delta \geq 0$, let the bias gradient of $\hat{t}$ satisfy the norm constraint $||\nabla b_{\underline{\theta}}||_C \leq \delta$, where $C$ is an arbitrary $n \times n$ symmetric positive definite matrix. Define $\mathcal{P}_B$ as the $n \times n$ matrix that projects onto the column space of $B = C^{-(1/2)} F_Y^+ C^{-(1/2)}$. Then, the variance $\sigma_{\underline{\theta}}^2$ of $\hat{t}$ satisfies

$$\sigma_{\underline{\theta}}^2 \geq B(\underline{\theta}, \delta) \quad (8)$$

where, if $\delta^2 \geq \nabla^T t_{\underline{\theta}} C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_{\underline{\theta}}$, then $B(\underline{\theta}, \delta) = 0$, whereas if $\delta^2 < \nabla^T t_{\underline{\theta}} C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_{\underline{\theta}}$ then

$$B(\underline{\theta}, \delta) = [\nabla t_{\underline{\theta}} + \underline{d}_{\min}]^T F_Y^+[\nabla t_{\underline{\theta}} + \underline{d}_{\min}],$$

$$= \lambda^2 \nabla^T t_{\underline{\theta}} C[\lambda C + F_Y^+]^{-1} F_Y^+[\lambda C + F_Y^+]^{-1} C \nabla t_{\underline{\theta}}, \quad (9)$$

where in (9)

$$\underline{d}_{\min} = -[\lambda C + F_Y^+]^{-1} F_Y^+ \nabla t_{\underline{\theta}}. \quad (10)$$
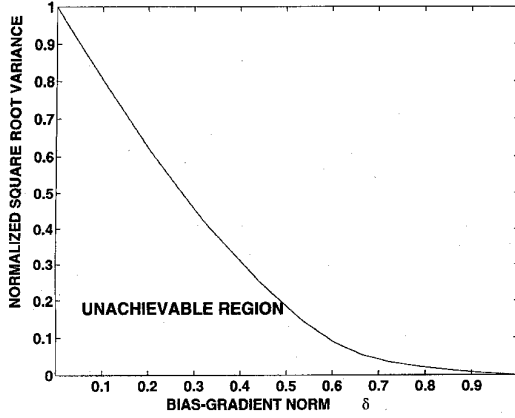
Fig. 1. Normalized uniform CR bound on the $\delta\sigma$ tradeoff plane for a specified value of $\underline{\theta}$.

In (9) and (10), $\lambda > 0$ is determined by the unique positive solution of $g(\lambda) = \delta^2$, where

$$g(\lambda) = \nabla^T t_{\underline{\theta}} F_Y^+ [\lambda C + F_Y^+]^{-1} C [\lambda C + F_Y^+]^{-1} F_Y^+ \nabla t_{\underline{\theta}}. \quad (11)$$

By tracing out the family of points $\{(\delta, \sqrt{B(\underline{\theta}, \delta)}): \delta \geq 0\}$, one obtains a curve in the $\delta\sigma$ plane for a particular $\underline{\theta} \in \Theta$. The curve is always monotone nonincreasing in $\delta$. Since $B(\underline{\theta}, \delta)$ is a lower bound on $\sigma_{\underline{\theta}}^2$, the region below the curve defines an unachievable region. Fig. 1 shows a typical delta-sigma tradeoff curve plotted in terms of normalized standard deviation $\sigma = \sqrt{B(\underline{\theta}, \delta)/B(\underline{\theta}, 0)}$. If an estimator lies on the curve, then lower variance can only be bought at the price of increased bias gradient and vice versa. For this reason, we call this curve the *delta-sigma tradeoff curve*.

It is important to point out that the delta-sigma tradeoff curve can be generated without solving the nonlinear equation $g(\lambda) = \delta^2$ (11), which typically must be solved via numerical methods. It is much easier to continuously vary $\lambda$ over the range $(0, \infty)$ and sweep out the curve by using the $\lambda$ parameterizations of $\delta^2$ and $B(\underline{\theta}, \delta)$ specified by relations (11) and (9), respectively.

*Comments:*

- The uniform bound $B(\underline{\theta}, \delta)$ is always less than or equal to the unbiased CR bound $B(\underline{\theta}; 0) = \nabla^T t_{\underline{\theta}} F_Y^+ \nabla t_{\underline{\theta}}$. The slope of $B(\underline{\theta}, \delta)$ at $\delta = 0$ gives a bias sensitivity index $\eta$ for the unbiased CR bound. For nonsingular $F_Y$ and single component estimation $(t_{\underline{\theta}} = \theta_1)$, it is shown in [15] that $\eta = 2\sqrt{1 + \underline{c}^T F_S^{-2} \underline{c}}$, where $\underline{c}$ is the first column of $F_Y$, and $F_S$ is the principal minor of the (1,1) element of $F_Y$. Large values of this index indicate that the unbiased form of the CR bound is not reliable for estimators that may have very small, and perhaps even unmeasurable, biases.

- The orthogonal projection $\mathcal{P}_B$ can be expressed either as $\mathcal{P}_B = B[B^T B]^+ B^T = B^+ B = B B^+$ or via the eigendecomposition of $B$ as $\mathcal{P}_B = \Sigma_{i=1}^r \underline{\xi}_i \underline{\xi}_i^T$, where $r$ is the rank of $F_Y$, and $\{\underline{\xi}_i\}_{i=1}^r$ are orthonormal eigenvectors associated with the nonzero eigenvalues of $B$. By using properties of the Moore–Penrose pseudo-inverse, it can

be shown that

$$\nabla^T t_{\underline{\theta}} C^{1/2} \mathcal{P}_B C^{1/2} \nabla t_{\underline{\theta}} = \nabla^T t_{\underline{\theta}} F_Y^+ [C^{-(1/2)} F_Y^+]^+ C^{1/2} \nabla t_{\underline{\theta}}.$$

- When $F_Y$ is nonsingular

$$F_Y^+ = F_Y^{-1},$$
$$\mathcal{P}_B = I, \nabla^T t_{\underline{\theta}} C^{1/2} \mathcal{P}_B C^{1/2} \quad \nabla t_{\underline{\theta}} = \|\nabla t_{\underline{\theta}}\|_C^2$$

and (9)–(11) of Theorem 1 reduce to

$$\begin{aligned}
B(\underline{\theta}, \delta) \\
&= [\nabla t_{\underline{\theta}} + \underline{d}_{\min}]^T F_Y^{-1} [\nabla t_{\underline{\theta}} + \underline{d}_{\min}] \\
&= \lambda^2 \nabla^T t_{\underline{\theta}} [C^{-1} + \lambda F_Y]^{-1} F_Y [C^{-1} + \lambda F_Y]^{-1} \nabla t_{\underline{\theta}}
\end{aligned}$$
$$(12)$$

where

$$\underline{d}_{\min} = -C^{-1} [C^{-1} + \lambda F_Y]^{-1} \nabla t_{\underline{\theta}} \quad (13)$$

and $\lambda > 0$ is given by the unique positive solution of $g(\lambda) = \delta^2$, where

$$g(\lambda) = \nabla^T t_{\underline{\theta}} [C^{-1} + \lambda F_Y]^{-1} C^{-1} [C^{-1} + \lambda F_Y]^{-1} \nabla t_{\underline{\theta}}. \quad (14)$$

When $C = I$ and $t_{\underline{\theta}} = \theta_1$, these are identical to the results obtained in [15].

- In Theorem 1, $\underline{d}_{\min}$ defined in (10) is an optimal bias gradient in the sense that it minimizes the biased CR bound (7) over all vectors $\nabla b_{\underline{\theta}}$, satisfying the constraint $\|\nabla b_{\underline{\theta}}\|_C \leq \delta$. The resultant bound is independent of the particular estimator bias as long as the bias gradient norm constraint holds. From the proof of Theorem 1, if $\delta^2 \geq \nabla^T t_{\underline{\theta}} C^{1/2} \mathcal{P}_B C^{1/2} \nabla t_{\underline{\theta}}$, then the minimizing bias gradient is of the form $\underline{d}_{\min} = -\mathcal{P}_B C^{1/2} \nabla t_{\underline{\theta}} + \phi$, where $\phi$ is any vector satisfying $B\phi = \underline{0}$, and $\|\phi\|^2 \leq \delta^2 - \nabla^T t_{\underline{\theta}} C^{1/2} \mathcal{P}_B C^{1/2} \nabla t_{\underline{\theta}}$. Thus, for the case of singular $F_Y$, there exist many optimal bias gradients.

- An estimator is said to locally achieve a bound in a neighborhood of a point $\underline{\theta}$ if the estimator achieves the bound whenever the true parameter lies in the neighborhood. It has been shown [15] that if $F_Y$ is nonsingular, if $\delta$ is small, if $t_{\underline{\theta}} = \theta_1$, and if the unbiased matrix CR bound is locally achievable by an unbiased estimator $\hat{\underline{\theta}}^*$ in a neighborhood of a point $\underline{\theta}$, then one can construct an estimator that locally achieves the uniform bound in this neighborhood by introducing a small amount of bias into $\hat{\underline{\theta}}^*$. However, since unbiased estimators may not exist for singular $F_Y$, the uniform CR bound for singular $F_Y$ may not be locally achievable. An example where the bound is globally achievable over all $\Theta$ is presented in Section IV.

- While we will not use it in this paper, a more general form of Theorem 1 holds for the case that $C$ may be nonnegative definite. This situation is relevant for cases where the user does not wish to penalize the estimator for high bias variation over certain hyperplanes in the parameter space. For example, when estimation of image contrast is of interest, spatially homogeneous biases may be tolerable, and $C$ may be chosen to be of rank $n - 1$ having the vector $\underline{1} = [1, \cdots, 1]^T$ in its nullspace. Let

$B(\underline{\theta}, \delta), \underline{d}_{\min}$ and $g(\lambda)$ be as defined in Theorem 1. Assume that $C$ is nonnegative definite, but $F_Y^+ + \lambda C$ is positive definite for $0 < \lambda < \infty$. For fixed $\delta > 0$, let the bias gradient of $\hat{t}$ satisfy the semi-norm constraint $\|\nabla b_{\underline{\theta}}\|_C \le \delta$. Then, $\mathrm{var}_{\underline{\theta}}(\hat{t}) \ge B^*(\underline{\theta}, \delta)$, where

$$B^*(\underline{\theta}, \delta) = \begin{cases} B(\underline{\theta}, g(\infty)), & 0 \le \delta^2 < g(\infty) \\ B(\underline{\theta}, \delta), & g(\infty) \le \delta^2 \le g(0) \\ 0, & g(0) < \delta^2 < \infty, \end{cases}$$

$g(\infty) = \lim_{\lambda \to \infty} g(\lambda)$, and $g(0) = \lim_{\lambda \to 0} g(\lambda)$.

### C. Recipes for Uniform Bound Computation

As written in Theorem 1, (9)–(11) are not in the most convenient form for computation as they involve several matrix multiplications and inversions. An equivalent form for the pair $B(\underline{\theta}, \delta)$ and $g(\lambda)$ in (9) and (11) was obtained in the process of proving the theorem ((47) and (48))

$$B(\underline{\theta}, \delta) = \lambda^2 \nabla^T t_{\underline{\theta}} C^{1/2} [\lambda I + B]^{-1} B [\lambda I + B]^{-1} C^{1/2} \nabla t_{\underline{\theta}}$$

$$g(\lambda) = \nabla^T t_{\underline{\theta}} [C^{1/2} B (\lambda I + B)^{-2} B C^{1/2}] \nabla t_{\underline{\theta}} = \delta^2$$

where $B = C^{-(1/2)} F_Y^+ C^{-(1/2)}$. If an eigendecomposition of the matrix $B$ is available, the delta-sigma tradeoff curve can be efficiently computed by sweeping out $\lambda$ in the following pair of weighted sums of inner products

$$B(\underline{\theta}, \delta) = \sum_{i=1}^{n} \frac{\lambda^2 \beta_i}{(\lambda + \beta_i)^2} |\nabla^T t_{\underline{\theta}} C^{1/2} \underline{\xi}_i|^2 \qquad (15)$$

$$\delta^2 = \sum_{i=1}^{n} \frac{\beta_i^2}{(\lambda + \beta_i)^2} |\nabla^T t_{\underline{\theta}} C^{1/2} \underline{\xi}_i|^2 \qquad (16)$$

where $\beta_i$ and $\underline{\xi}_i$ denote an eigenvalue and eigenvector of $B$.

When $F_Y$ is positive definite but ill conditioned, the computation of $B$ may be numerically unstable. In this case, it is better to use the equivalent form

$$B(\underline{\theta}, \delta) = \lambda^2 \nabla^T t_{\underline{\theta}} C^{1/2} [I + \lambda G]^{-1} G [I + \lambda G]^{-1} C^{1/2} \nabla t_{\underline{\theta}} \tag{17}$$

$$\delta^2 = \nabla^T t_{\underline{\theta}} C^{1/2} [I + \lambda = G]^{-2} C^{1/2} \nabla t_{\underline{\theta}} \tag{18}$$

where $G = B^{-1} = C^{1/2} F_Y C^{1/2}$. Note that computation of the form (17) and (18) requires only a single matrix inversion $[I + \lambda G]^{-1}$. Since $\lambda > 0$ and $F_Y$ is positive definite, this inversion is well conditioned except possibly if $\lambda$ is very large.

The eigendecomposition of $G$ can be used in (17) and (18) to produce a pair of expressions similar to (15) and (16) for computing the delta-sigma tradeoff curve for positive definite $F_Y$. Alternatively, the right-hand sides of (17) and (18) can be approximated by using iterative equation solving methods such as Gauss–Seidel (GS) or preconditioned conjugate gradient (CG) algorithms [35]. See [45] and [36] for a more detailed discussion of the application of iterative equation solvers to CR bound approximation. This approach can be implemented in the following sequence of steps.

1) Select $\lambda \in (0, \infty)$.
2) Compute $\underline{x} = [I + \lambda F_Y C]^{-1} \nabla t_{\underline{\theta}}$ by applying CG or GS iterations to solve the following linear equation for $\underline{x}$:

$$[I + \lambda F_Y C] \underline{x} = \nabla t_{\theta}.$$

3) Compute $\underline{y} = C \underline{x}$.
4) Compute the point $(\delta, B(\underline{\theta}, \delta))$ via

$$B(\underline{\theta}, \delta) = \lambda^2 \underline{x}^T F_Y \underline{x}$$

$$\delta = \sqrt{\underline{x}^T \underline{y}}.$$

Since step 2 must be repeated for each value of $\lambda$, this method is competitive only when one is interested in evaluation of the curve $B(\underline{\theta}, \delta)$ at a small number of values of $\delta = \delta(\lambda)$. When a denser sampling of the curve is desired an eigendecomposition method, e.g., as in (15) and (16), becomes more attractive since once the quantities $\beta_i$ and $|\nabla^T t_{\underline{\theta}} C^{1/2} \underline{\xi}_i|^2$ are available, the curve can be swept out over $\lambda$ without performing additional vector operations.

### III. ESTIMATION OF BIAS GRADIENT NORM

To be able to compare the performance of an estimator against the uniform CR bound of Theorem 1, we need to determine the estimator variance and the bias gradient length. In most cases, the bias gradient cannot be determined analytically, and it is therefore important to have a computationally efficient method to estimate it either experimentally or via simulations. A brute force estimate would be to estimate the finite difference approximation

$$\nabla b_{\underline{\theta}} \approx \frac{1}{\epsilon} [b_{\underline{\theta} + \epsilon \underline{e}_1} - b_{\underline{\theta}}, \cdots, b_{\underline{\theta} + \epsilon \underline{e}_n} - b_{\underline{\theta}}]$$

but this requires performing a seperate simulation run for each coordinate perturbation $\underline{\theta} + \epsilon \underline{e}_k$. In the following, we describe a more direct method for estimating the bias gradient that does not require performing multiple simulation runs nor does it require making a finite difference approximation. The method is based on the fact that for any random variable $Z$ with finite mean

$$E_{\underline{\theta}} \left[ Z \frac{\partial \ln f_Z(Z; \underline{\theta})}{\partial \theta_j} \right] = \int_Z z \frac{\partial f_Z(z; \underline{\theta})}{\partial \theta_j} \, dz$$

$$= \frac{\partial}{\partial \theta_j} \int_Z z f_Z(z; \underline{\theta}) \, dz$$

$$= \frac{\partial}{\partial \theta_j} E_{\underline{\theta}}(Z). \qquad (19)$$

Thus, in particular, we have the following relation:

$$\nabla b_{\underline{\theta}} = E_{\underline{\theta}}[\hat{t}(Y) \nabla_{\underline{\theta}} \ln f_Y(Y; \underline{\theta})] - \nabla t_{\underline{\theta}}.$$

Since $E_{\underline{\theta}}[\nabla_{\underline{\theta}} \ln f_Y(Y; \underline{\theta})] = 0$, an equivalent relation is

$$\nabla b_{\underline{\theta}} = E_{\underline{\theta}}[(\hat{t}(Y) - \zeta) \nabla_{\underline{\theta}} \ln f_Y(Y; \underline{\theta})] - \nabla t_{\underline{\theta}} \qquad (20)$$

for any random variable $\zeta$ statistically independent of $Y$. As explained in the following discussion, the quantity $\zeta$ can be used to control the variance of the bias gradient estimate.

Substituting sample averages for ensemble averages in (20), we obtain the following unbiased and consistent estimator of the bias gradient vector $\nabla b_{\underline{\theta}}$

$$\widehat{\nabla b_{\underline{\theta}}} = \frac{1}{L} \sum_{i=1}^{L} (\hat{t}(Y_i) - \zeta_i) \nabla_{\underline{\theta}} \ln f_Y(Y_i; \underline{\theta}) - \nabla t_{\underline{\theta}} \qquad (21)$$

where $\{Y_i\}_{i=1}^{L}$ is a set of i.i.d. realizations from $f_Y(y;\underline{\theta})$. In (21), $\{\zeta_i\}_{i=1}^{L}$ is any sequence of i.i.d. random variables such that $Y_i, \zeta_i$ are statistically independent for each $i$.

It can be shown that when $\zeta_i = 0$ for all $i$ the covariance matrix of $\widehat{\nabla b_{\underline{\theta}}}$ is the matrix sum

$$S(\widehat{\nabla b_{\underline{\theta}}}) = \frac{1}{L} \text{cov}_{\underline{\theta}}\left((\hat{t}(Y_i) - m_{\underline{\theta}})\nabla \ln f_{Y_i}(Y_i;\underline{\theta})\right)$$
$$+ \frac{1}{L}[2m_{\underline{\theta}}\boldsymbol{R} + m_{\underline{\theta}}^2 \boldsymbol{F}_{Y_i}] \qquad (22)$$

where

$$\boldsymbol{R} = E_{\underline{\theta}}[(\hat{t}(Y_i) - m_{\underline{\theta}})[\nabla_{\underline{\theta}} \ln f_{Y_i}(Y_i;\underline{\theta})][\nabla_{\underline{\theta}}^T \ln f_{Y_i}(Y_i;\underline{\theta})]^T]$$

and

$$\boldsymbol{F}_{Y_i} = E_{\underline{\theta}}([\nabla_{\underline{\theta}} \ln f_Y(Y_i;\underline{\theta})][\nabla_{\underline{\theta}} \ln f_Y(Y_i;\underline{\theta})]^T)$$

is the single trial Fisher information. The first term on the right-hand side (RHS) of (22) decays as $1/L$ and is independent of the mean $m_{\underline{\theta}}$. The second term also decays as $1/L$ but is unbounded in the mean $m_{\underline{\theta}}$. It is easily shown that this term can be eliminated by setting $\zeta_i = m_{\underline{\theta}} = $ constant in (21) but this is not a practical since the mean $m_{\underline{\theta}} = E_{\underline{\theta}}[\hat{t}(Y_i)]$ is unknown to the user. However, we can use the punctured sample mean estimate:

$$\zeta_i = \frac{1}{L-1}\sum_{\substack{j=1 \\ j \neq i}}^{L} \hat{t}(Y_j)$$

which is unbiased and, as required for the validity of (20), is statistically independent of $Y_i$. Substitution of this $\zeta_i$ into (21) gives, after simplification, the following unbiased and consistent sample mean estimate of $\nabla b_{\underline{\theta}}$:

$$\widehat{\nabla b_{\underline{\theta}}} = \frac{1}{L-1}\sum_{i=1}^{L}\left(\hat{t}(Y_i) - \frac{1}{L}\sum_{j=1}^{L}\hat{t}(Y_j)\right)\nabla_{\underline{\theta}} \ln f_Y(Y_i;\underline{\theta})$$
$$- \nabla t_{\underline{\theta}}. \qquad (23)$$

A simple calculation shows that the covariance of (23) is

$$S(\widehat{\nabla b_{\underline{\theta}}}) = \frac{1}{L}\text{cov}_{\underline{\theta}}\left((\hat{t}(Y_i) - m_{\underline{\theta}})\nabla_{\underline{\theta}} \ln f_Y(Y_i;\underline{\theta})\right)$$
$$+ \frac{1}{L(L-1)}[\text{var}_{\underline{\theta}}(\hat{t}(Y_i))\boldsymbol{F}_{Y_i} + \nabla b_{\underline{\theta}}\nabla^T b_{\underline{\theta}}]. \qquad (24)$$

Note that the second term in (24) depends on $m_{\underline{\theta}}$ only through its gradient and decreases to zero at the much faster asymptotic rate of $1/L^2$ as compared with the rate $1/L$ in (22).

### A. A Bootstrap Estimator for Bias Gradient Norm

A natural "method-of-moments" estimate for $\delta_{\theta}^2 = ||\nabla b_{\underline{\theta}}||_C^2$ is the norm squared of the unbiased estimator $\hat{\delta}^2 = ||\widehat{\nabla b_{\underline{\theta}}}||_C^2$ (21). It can easily be shown that this estimator is biased with

bias equal to $E_{\underline{\theta}}[||\widehat{\nabla b_{\underline{\theta}}} - \nabla b_{\underline{\theta}}||_C^2] = \text{trace}\{S(\widehat{\nabla b_{\underline{\theta}}})\}$, which, in view of (22) or (24), decays to zero only as $1/L$. Below, we present a norm estimator based on the bootstrap resampling methodology whose bias decays at a faster rate.

Let $Y_1^*, \cdots, Y_L^*$ denote a *bootstrap sample* obtained by randomly resampling the realizations $Y_1 = y_1, \cdots, Y_L = y_L$ with replacement. Given the estimate $\hat{\delta}^2 = \hat{\delta}^2(y_1, \cdots, y_L) = ||\widehat{\nabla b_{\underline{\theta}}}||_C^2$, the bootstrap estimate of $\delta_{\theta}^2$ is defined as the expectation of $\hat{\delta}_*^2 = \hat{\delta}^2(Y_1^*, \cdots, Y_L^*)$ with respect to the resampling distribution [37]

$$E_*[\hat{\delta}_*^2] = \sum_* \hat{\delta}^2(Y_1^*, \cdots, Y_L^*)\binom{L}{c_1 \cdots c_L}L^{-L}. \qquad (25)$$

In (25), $c_i$ is the number of times the value $y_i$ appears in the set $\{Y_j^*\}_{j=1}^{L}$, and $\Sigma_*$ denotes a summation over all nonnegative integers $c_1, \cdots, c_L$ satisfying $\Sigma_{i=1}^{L}c_i = L$. The bootstrap estimate of the bias of the estimator $\hat{\delta}^2$ is defined as $E_*[\hat{\delta}_*^2] - \hat{\delta}^2$, which leads to the bias corrected estimator $\hat{\delta}_c^2$

$$\hat{\delta}_c^2 = 2\hat{\delta}^2 - E_*[\hat{\delta}_*^2]. \qquad (26)$$

Due to the simple quadratic dependence of $\hat{\delta}^2$ on the single sample quantities $\hat{t}(y_i)\nabla_{\underline{\theta}} \ln f_Y(y_i;\underline{\theta}), i = 1, \cdots, L$, the expectation (25) can be expressed in analytical form (see Appendix B), leading to the bias-corrected estimate

$$\hat{\delta}_c^2 = ||\widehat{\nabla b_{\underline{\theta}}}||_C^2 - \frac{1}{L^2}\sum_{i=1}^{L}||\widehat{\nabla b_{\underline{\theta}}}(y_i) - \widehat{\nabla b_{\underline{\theta}}}||_C^2 \qquad (27)$$

where $\widehat{\nabla b_{\underline{\theta}}}(y_i)$ is the estimate (21) based on a single sample $(\zeta_i = 0): \hat{t}(y_i)\nabla_{\underline{\theta}} \ln f_Y(y_i,\underline{\theta}) - t_{\underline{\theta}}$. The bias of $\hat{\delta}_c^2$ is equal to

$$E_{\underline{\theta}}[\hat{\delta}_c^2] - ||\nabla b_{\underline{\theta}}||_C^2 = \frac{1}{L}E_{\underline{\theta}}[||\widehat{\nabla b_{\underline{\theta}}} - \nabla b_{\underline{\theta}}||_C^2]$$

which, relative to the estimator $||\widehat{\nabla b_{\underline{\theta}}}||_C^2$, decays to zero at the much faster rate of $1/L^2$. However, if $L$ is insufficiently large, the bootstrap estimator $\hat{\delta}_c^2$ may take on negative values.

## IV. APPLICATION TO INVERSE PROBLEMS

We use the theory developed above to perform a study of fundamental bias-variance tradeoffs for three general classes of inverse problems. First, we consider well-posed linear Gaussian inverse problems that have positive definite Fisher information. Next, we consider ill-posed Gaussian inverse problems where the Fisher matrix is singular. For these two linear applications, an exact analysis is possible since all curves in the delta-sigma tradeoff plane have analytic expressions. Finally, we study a nonlinear Poisson inverse problem to illustrate the empirical bias-gradient norm approximations discussed in the previous section.

### A. Linear Gaussian Model

Assume that the observation consists of a vector $Y = \underline{Y} \in \mathbb{R}^m$ that obeys the linear Gaussian model:

$$\underline{Y} = \boldsymbol{A}\underline{\theta} + \underline{\epsilon}, \qquad (28)$$

where,

**A**  $m \times n$ coefficient matrix called the system matrix;

$\underline{\theta}$  unknown source;

$\underline{\epsilon}$  vector of zero mean Gaussian random variables with positive definite covariance matrix $\Sigma$.

For concreteness, we will refer to $\theta_i$ as the intensity of the source at pixel $i$. The Fisher information matrix has the well-known form [19]

$$F_Y = A^T \Sigma^{-1} A. \tag{29}$$

This matrix is nonsingular when $A$ is of full column rank $n$. We will consider estimation of the linear combination $t_\theta = \underline{h}^T \underline{\theta}$, where $\underline{h}$ is a fixed nonzero vector in $\mathbb{R}^n$. Since $F_Y$ and $\nabla t_\theta = \underline{h}$ are not functionally dependent on $\underline{\theta}$, the uniform bound $B(\underline{\theta}, \delta)$ will not depend on the specific form of the unknown source $\underline{\theta}$.

To demonstrate the achievability of the fundamental delta-sigma tradeoff curve, we consider the quadratically penalized maximum likelihood (QPML) estimator. The QPML strategy is frequently used in order to obtain stable solutions in the presence of small variations in experimental conditions [38] or as a way to incorporate parameter constraints or *a priori* information [39]. For the linear Gaussian problem (28), the QPML estimator of the linear combination $t_{\underline{\theta}} = \underline{h}^T \underline{\theta}$ is $\hat{t} = \underline{h}^T \hat{\underline{\theta}}$, where $\hat{\underline{\theta}}$ minimizes the following objective function over $\underline{\theta}$:

$$[\underline{Y} - A\underline{\theta}]^T \Sigma^{-1} [\underline{Y} - A\underline{\theta}] + \beta \underline{\theta}^T P \underline{\theta}. \tag{30}$$

In the above, $\beta > 0$ is a regularization parameter, and $P$ is a symmetric nonnegative-definite penalty matrix. For ill-conditioned or singular $A$, the penalty improves the numerical stability of the matrix inversion $[F_Y + \beta P]^{-1}$ in (31) below by lowering its condition number. The simplest choice for the penalty matrix $P$ is the identity $I$, which yields a class of energy penalized least squares estimators variously known as Tikonov regularized least squares in the inverse problem literature [38] and shrinkage estimation or ridge regression in the multivariate statistics literature [5]. A popular choice in imaging applications is to use a nondiagonal differencing type operator to enforce smoothness constraints [40], [41].

The minimizer of (30) is the penalized weighted least squares (PLS) estimator

$$\hat{\underline{\theta}} = [F_Y + \beta P]^{-1} A^T \Sigma^{-1} \underline{Y} \tag{31}$$

yielding the QPML estimator $\hat{t} = \underline{h}^T \hat{\underline{\theta}}$.

The estimator bias is

$$b_{\underline{\theta}} = \underline{h}^T [[F_Y + \beta P]^{-1} F_Y - I]\underline{\theta}$$

and its bias gradient is

$$\nabla b_{\underline{\theta}} = [F_Y [F_Y + \beta P]^{-1} - I]\underline{h} \tag{32}$$

$$= -\beta P [\beta P + F_Y]^{-1} \underline{h}. \tag{33}$$

Finally, the variance of the QPML estimator $\hat{t}$ is

$$\sigma_{\underline{\theta}}^2 = \frac{1}{\beta^2} \underline{h}^T \left[ P + \frac{1}{\beta} F_Y \right]^{-1} F_Y \left[ P + \frac{1}{\beta} F_Y \right]^{-1} \underline{h}. \tag{34}$$

Consider the special case of estimation of a single component $\theta_k$ of $\underline{\theta}$ for which $\underline{h} = \underline{e}_k = [0, \cdots, 0, 1, 0, \cdots, 0]^T$. When the matrices $F_Y$ and $P$ commute, as occurs, for example, when $P = I$, the bias gradient (32) is seen to be equal to the difference between the mean response $[\beta P + F_Y]^{-1} F_Y \underline{e}_k$ of the PLS estimator to a point source $\underline{\theta} = \underline{e}_k$, i.e., the point spread function of the estimator and the ideal point response $\underline{e}_k$. Thus, under the commutative assumption, the bias gradient norm can be viewed as a measure of the geometric resolution of $\hat{t}$ [16].

*1) Positive Definite Fisher Matrix:* Assume that $F_Y$ is positive definite, and compare (33) and (34) to (13) and (12) for $\underline{d}_{\min}$ and the bound $B(\underline{\theta}, \delta)$, respectively. Identifying $\nabla t_\theta = \underline{h}, \lambda = 1/\beta$, it is clear that when $P$ is chosen as $C^{-1}$, the PLS estimator achieves the bound $B(\underline{\theta}, \delta)$ and has optimal bias gradient $\underline{d}_{\min}$. Thus, for linear functions $t_\theta$, the uniform bound is achievable, and the region above and including the fundamental delta-sigma tradeoff curve is an *achievable region*. Furthermore, since the bias gradient is a linear function, from (6), we have a very strong optimality property: The QPML estimator $\hat{t}$ is a minimum variance *biased* estimator in the sense that it is an estimator of minimum variance among estimators that satisfy the maximal bias constraint $\sup_{\underline{u} \in \mathcal{C}} |b_{\underline{u}} - b_{\underline{\theta}}| \leq \delta$, where $\delta^2 = g(1/\beta)$ and $\mathcal{C}$ is the ellipsoid defined above (3).

We used the computational recipe presented in Section II to trace out the delta-sigma tradeoff curve (uniform bound) parametrically as a function of $\lambda > 0$. Fig. 2 shows the delta-sigma tradeoff curve for the case of pixel intensity estimation ($\underline{h} = \underline{e}_{67}, \Sigma = I$) and a well-conditioned full rank discrete Gaussian system matrix. Specifically, we generated a $128 \times 128$ matrix $A$ with elements $a_{ij} = (1/\sqrt{2\pi}w)e^{-(i-j)^2/2w^2}$ and $w = 0.5$. The condition number of $A$ is 1.7. The matrix $C$ in the norm $||\nabla b_{\underline{\theta}}||_C$ was selected as the inverse of the second-order (Laplacian) differencing matrix

$$C^{-1} = \begin{bmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & 0 & \\ & \ddots & \ddots & \ddots & \\ 0 & -1 & 2 & -1 \\ & & -1 & 2 \end{bmatrix}. \tag{35}$$

With this norm, the restriction $||\nabla b_{\underline{\theta}}||_C \leq \delta$ corresponds to a constraint on maximal bias variation $\max_{\theta \in \mathcal{C}} |\Delta b_{\underline{\theta}}|$ over a roughness-constrained neighborhood $\mathcal{C}(\underline{\theta}, \overline{C})$ of $\underline{\theta}$ (recall relation (6)). The performance curves $(||\nabla b_{\underline{\theta}}||_C, \sigma_{\underline{\theta}})$ for two PLS pixel intensity estimators (31) (one using the smoothing matrix $P = C^{-1}$, called the *smoothed QPML estimator* and another using the diagonal "energy penalty" $P = I$ called the *unsmoothed QPML estimator*) are also plotted in Fig. 2. These curves were traced out in the bias variance tradeoff plane by varying $\beta$ in the parametric descriptions of estimator variance (34) and estimator bias gradient (33).

*2) Singular Fisher Matrix:* When $A$ has rank less than $n$, $F_Y$ is singular, and unbiased estimators may not exist for all linear functions $t_{\underline{\theta}}$ of $\underline{\theta}$ [19], [42]. A lower bound on the norm of the bias gradient can be derived (see Appendix C) using the relation (6) between the norm and the maximal bias variation over a region of parameter space. Since the uniform
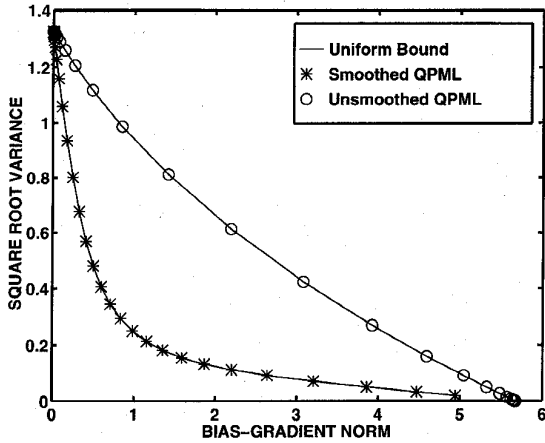
Fig. 2. Bias variance tradeoff study for pixel-intensity estimation and non-singular Fisher information. The smoothed PLS estimator (labeled smoothed QPML) exactly achieves the uniform bound.
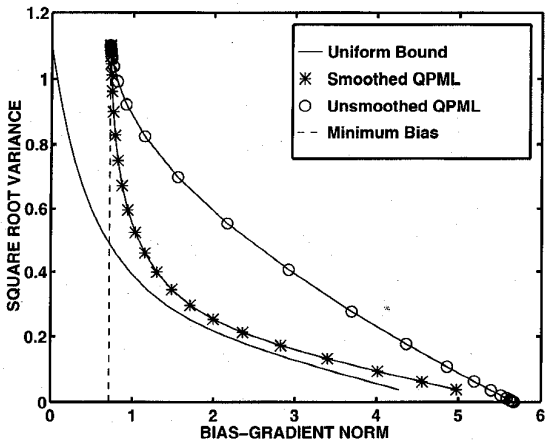


Fig. 3. Bias variance tradeoff study for pixel intensity estimation and singular Fisher information. Neither of the QPML estimators achieve the uniform bound.



Fig. 4. Coefficients $\underline{h}$ of contrast function $t_{\underline{\theta}} = \underline{h}^T \underline{\theta}$ used for bias variance tradeoff study shown in Fig. 5.



Fig. 5. Bias variance tradeoff study for estimation of contrast function (illustrated in Fig. 4) for singular Fisher information. The smoothed QPML estimator of contrast virtually achieves the uniform bound for bias-gradient norm $\delta$ greater than 0.2.

CR bound is finite and equal to the unbiased CR bound at $\delta = 0$, we cannot expect the delta-sigma tradeoff curve to be achievable for all $\delta$ as in the nonsingular case.

To illustrate we repeat the study of Fig. 2 with a rank deficient Gaussian kernel matrix $A_r$ obtained by decimating the rows of a full-rank Gaussian kernel matrix $A$ ($w = 2$) by a factor of 4. This yields the ill-posed problem of estimating a vector of 128 pixel intensities $\underline{\theta}$ based on only 32 observations $\underline{Y}$. We used the singular value decomposition of $A$ to compute the delta-sigma tradeoff curve and the minimal bias gradient norm. The results for pixel intensity estimation ($t_{\underline{\theta}} = \underline{e}_{67}$) are plotted in Fig. 3 along with the performance curves associated with smoothed QPML ($P = C^{-1}$ of (35)) and unsmoothed QPML ($P = I$) estimators. Note that neither of the estimators achieve the uniform bound for any value of the parameter $\beta$. The bound on bias gradient norm (dashed line) is an asymptote on estimator performance that forces a sharp knee in the estimator performance curves. At points close to this knee, maximal reduction in bias is only achieved at the price of significant increase in the variance.
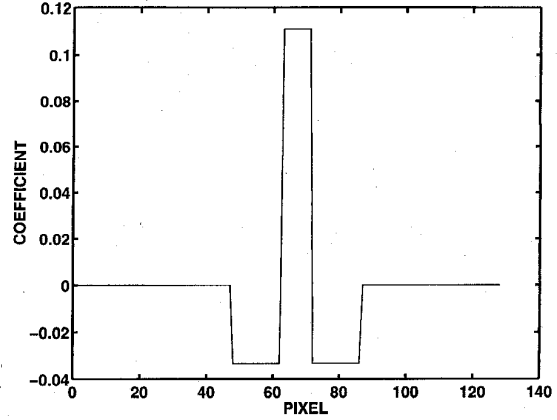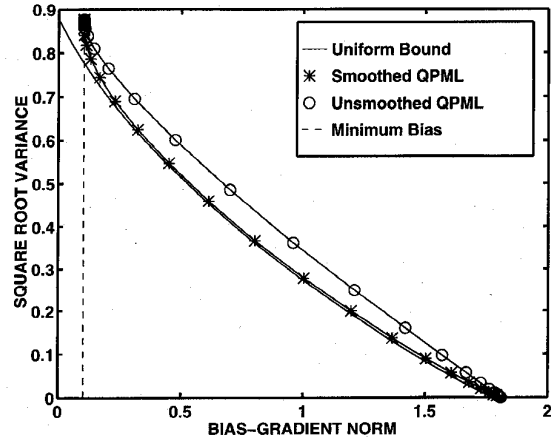
For comparison, in Fig. 5, we plot the analogous curves for smoothed and unsmoothed QPML estimation of the contrast function defined as $t_{\underline{\theta}} = \underline{h}^T \underline{\theta}$, where the elements of $\underline{h}$ are plotted in Fig. 4. Observe that the smoothed QPML estimator of contrast comes much closer to the uniform bound than does the smoothed QPML estimator of pixel intensity shown in Fig. 3.

Under certain conditions, the uniform CR bound is exactly achievable even for singular $F_Y$, although generally not by a QPML estimator of the form (31) and generally not for all $\delta$. Consider the estimator

$$\hat{\underline{\theta}} = \underline{h}^T [I + \beta F_Y^+ P]^{-1} F_Y^+ A^T \Sigma^{-1} \underline{y}. \qquad (36)$$

This estimator reduces to the previous estimator (31) for the case of nonsingular $F_Y$. The estimator bias gradient is

$$\nabla b_{\underline{\theta}} = (F_Y F_Y^+ [I + \beta P F_Y^+]^{-1} - I) \nabla t_{\underline{\theta}}$$

$$= -\left[\frac{1}{\beta} P^{-1} + F_Y^+\right]^{-1} F_Y^+ \underline{h}$$

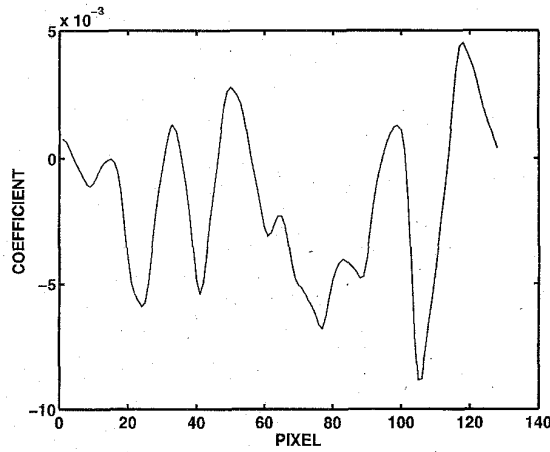$$- [I - F_Y F_Y^+][I + \beta P F_Y^+]^{-1} \underline{h} \qquad (37)$$

Fig. 6. Coefficients $\underline{h}$ of the linear compound $t_{\underline{\theta}} = \underline{h}^T \underline{\theta}$ satisfying condition 2 of Theorem 2 and used for computing the curves in Fig. 7.

and the estimator variance is

$$\sigma_{\underline{\theta}}^2 = \underline{h}^T P^{-1} [P^{-1} + \beta F_Y^+]^{-1} F_Y^+ [P^{-1} + \beta F_Y^+]^{-1} P^{-1} \underline{h} \quad (38)$$

where in (38), we have used the property $F_Y^+ F_Y F_Y^+ = F_Y^+$ [43]. Noting that here $\nabla t_{\underline{\theta}} = \underline{h}$, we conclude that the estimator variance is equal to the lower bound expression $B(\underline{\theta}, \delta)$ given in (9) when $P = C^{-1}$ and $\beta = 1/\lambda$. Furthermore, under these conditions, the bias gradient (37) differs from the optimal bias gradient $\underline{d}_{\min}^T$, which is given in (10), only by the presence of the second additive term on the right-hand side of (37). Thus, the estimator (36) with $P = C^{-1}$ is an optimal biased estimator when this second additive term is equal to zero.

We summarize these results in a theorem that applies to both singular and nonsingular $F_Y$.

*Theorem 2:* Let $B = P^{1/2} F_Y^+ P^{1/2}$, where $F_Y$ is the possibly singular Fisher information matrix. If

1) $\delta^2 < \underline{h}^T P^{1/2} \mathcal{P}_B P^{1/2} \underline{h}$ and
2) the vector $\underline{h}$ lies in the nullspace of $[I - F_Y F_Y^+][I + \beta P F_Y^+]^{-1}$

then the estimator $\hat{t}$ of $t_{\underline{\theta}} = \underline{h}^T \underline{\theta}$ given by (36) achieves the fundamental delta-sigma tradeoff in the sense of having minimum variance over all estimators satisfying $\|\nabla b_{\underline{\theta}}\|_{P^{-1}}^2 \leq \delta^2 = g(1/\beta)$, where $g(\cdot)$ is the function given in (11).

Recognizing the matrix $I - F_Y^+ F_Y = I - F_Y F_Y^+$ as the operator that projects onto the null space of $F_Y$, an equivalent condition to (2) is that $[I + \beta P F_Y^+]^{-1} \underline{h}$ lie in the range space of $F_Y$. For the special case of $\beta = 0$, condition (2) of Theorem 2 reduces to the well-known necessary condition for achievability of the unbiased CR bound: The vector $\underline{h}$ must lie in the range space of the Fisher information $F_Y$. In order for the uniform CR bound $B(\underline{\theta}, \delta)$ to be achievable for all values of $\delta$, condition (2) must hold for all $\beta > 0$. This is a much stronger condition except when the nullspace is independent of $\beta$, as occurs when $P = I$. This suggests that when estimation of any fixed $t_{\underline{\theta}}$ is of interest and the Fisher information is singular, the uniform bound will rarely be achievable everywhere in the delta-sigma plane.

To illustrate Theorem 2, we selected a small value of $\beta$ and found a vector $\underline{h}$ lying in the nullspace of the matrix $[I - F_Y^+ F_Y][I + \beta P F_Y^+]^{-1}$ via singular value decomposition.
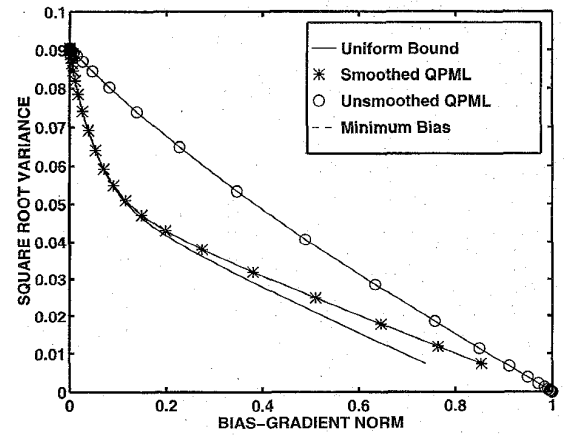


Fig. 7. Bias variance tradeoff study for estimation of nullspace compound $\underline{h}^T \underline{\theta}$ for singular Fisher information ($\underline{h}$ is illustrated in Fig. 6). The bound is exactly achieved by the smoothed QPML estimator at the point $\delta = 5.4 \times 10^{-6}$.

This vector is shown in Fig. 6. In view of Theorem 2, we know that the estimator (36) of $\underline{h}^T \underline{\theta}$ should achieve the uniform bound for the chosen value of $\beta$. In Fig. 7, we plot the uniform bound for estimators of $\underline{h}^T \underline{\theta}$ and the performance curve of smoothed ($P = C^{-1}$) and unsmoothed ($P = I$) estimators of the form (36). Observe that the smoothed estimator essentially achieves the uniform bound for $\delta < 0.2$.

### B. Poisson Model

In some applications, the observations $\underline{Y}$ are given by the linear model (28) but with non-Gaussian additive noise. Here, we consider the case of Poisson noise that arises in emission-computed tomography and other quantum-limited inverse problems [44]. The observation $\underline{Y} = [Y_1, \cdots, Y_m]^T$ is a vector of integers or counts with a vector of means $\underline{\mu} = [\mu_1, \cdots, \mu_m]^T$. This vector of counts obeys independent Poisson statistics with log-likelihood

$$\ln f_{\underline{Y}}(\underline{y}; \underline{\theta}) = \sum_{j=1}^{m} \left[ y_j \ln (\mu_j(\underline{\theta})) - \mu_j(\underline{\theta}) \right] + c. \quad (39)$$

In (39), $c$ is a constant independent of the unknown source $\underline{\theta}$, and the mean number of counts is assumed to obey the linear model

$$\underline{\mu}(\underline{\theta}) = A\underline{\theta} + \underline{r}. \quad (40)$$

For example, in emission-computed tomography

$\underline{\mu}$    vector of mean object projections measured over $m$ detectors;

$A$    $m \times n$ system matrix that depends on the tomographic geometry

$\underline{\theta}$    unknown image intensity vector;

$\underline{r}$    $m \times 1$ vector representing background noise due to randoms and scattered photons.

The Fisher information has the form [45]

$$F_Y(\underline{\theta}) = \sum_{j=1}^{m} \frac{1}{\mu_j(\underline{\theta})} A_{j*} A_{j*}^T, \quad (41)$$

where $A_{j*}^T$ is the $j$th row of $A$.

To investigate the achievability of the region above the delta–sigma tradeoff curve and to illustrate the empirical computation of bias gradient, we consider again the QPML strategy. The QPML estimator studied is $\hat{t} = t_{\hat{\theta}}$, where $\hat{\theta}$ is the vector $\underline{\theta}$ that maximizes the penalized likelihood function

$$J(\underline{\theta}) = \ln f_{\underline{Y}}(y; \underline{\theta}) - \frac{\beta}{2}\underline{\theta}^T P\underline{\theta} \qquad (42)$$

where $P$ is a nonnegative definite matrix.

Exact analytic expressions for the variance, bias, and bias gradient of the QPML estimator are intractable. However, it will be instructive to consider asymptotic approximations to these quantities. In Appendix D, expressions for asymptotic bias, bias gradient, and variance are derived under the assumption that the difference between the projection $AE_\theta[\hat{\theta}]$ of the mean QPML image and the projection $A\underline{\theta}$ of the true image is small—frequently a very good approximation in image restoration and tomography. Specializing the results (58)–(60) in Appendix D to the case of linear functions $t_\theta = \underline{h}^T\underline{\theta}$, we obtain the following expressions for the asymptotic variance of $\hat{t}$:

$$\sigma_{\underline{\theta}}^2 = \underline{h}^T[F_Y(\underline{\theta}) + \beta P]^{-1}F_Y(\underline{\theta})[F_Y(\underline{\theta}) + \beta P]^{-1}\underline{h} \qquad (43)$$

and the asymptotic bias gradient

$$\nabla b_{\underline{\theta}} = -P\left[P + \frac{1}{\beta}F_Y(\underline{\theta})\right]^{-1}\underline{h} + O\left(\frac{1}{\beta}\right) \qquad (44)$$

where $O(1/\beta)$ is a remainder term of order $1/\beta$.

When we identify $P = C^{-1}$ and $\beta = 1/\lambda$, we see that the estimator variance is identical to the optimal variance (12) and that for linear $t_\theta$, the bias gradient is identical to the optimal bias gradient (13) to order $O(1/\beta)$. Therefore, assuming the bias gradient and variance approximations (44) and (43) are accurate, for linear $t_\theta$, we can expect that the fundamental delta-sigma tradeoff curve will be approximately achieved by the QPML estimator for large values of the regularization parameter $\beta$ if $P = C^{-1}$.

To examine the performance of the methods for estimating bias gradient norm described in Section III and to verify the asymptotic bias and variance performance predictions, we generated simulated Poisson measurements with means given by (40). In these simulations, $A$ was a $128 \times 128$ tridiagonal blurring matrix with kernel $(0.23, 0.54, 0.23)$ for which the condition number is 12.5. The source intensity $\underline{\theta}$ is shown in Fig. 8. The function of interest was chosen as $t_\theta = \theta_{65}$, which is the intensity of pixel 65 in Fig. 8. We generated $L = 1000$ realizations of the measurements, each having a mean total of $\Sigma_{j=1}^m \mu_j(\underline{\theta}) = 2100$ counts, including a 5% background representing random coincidences [20].

We computed three types of estimates of $\underline{\theta}$:

i) the quadratically penalized maximum likelihood estimator using the "energy penalty" $(P = I)$;

ii) a truncated SVD estimator;

iii) a "deconvolve/shrink" estimator.



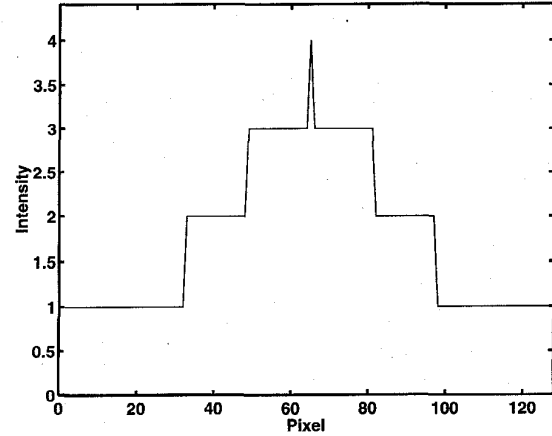Fig. 8. Emission source $\underline{\theta}$ used for Poisson simulations. The spike in the center was the pixel of interest.

We maximized the nonquadratic penalized likelihood objective using the PML-SAGE algorithm, which is a variant of the iterative space alternating generalized expectation-maximization (SAGE) algorithm of [20] adapted for penalized maximum-likelihood image reconstruction [46]. We initialized PML-SAGE with an unweighted penalized least-squares estimate: $(A^TA + \beta^*I)^{-1}A'(Y - \underline{r})$, which is linear so that it can be computed noniteratively. Here, $\beta^* = \beta \Sigma_j A_{jk}/(\Sigma_j A_{jk}/Y_j)$ for $k = 65$ (cf. [47] and [48]). By so initializing, only 30 iterations were needed to ensure convergence to a precision well below the estimate standard deviation. For the truncated singular value decomposition estimator, we computed the singular value decomposition of $A$ and computed the approximate pseudoinverse of $A$ by excluding the 10 smallest eigenvalues. The form of the "deconvolve/shrink" estimator is

$$\hat{\underline{\theta}}(Y) = \beta(A^TA)^{-1}A^T(Y - \underline{r})$$

where $\beta$ ranges from 0 to 1.

We applied each estimator to the $L = 1000$ measurement realizations and computed the standard sample variance

$$\hat{\sigma}^2 = \frac{1}{L-1}\sum_{i=1}^{L}(\hat{t}(Y_i) - \bar{\hat{t}})^2$$

where

$$\bar{\hat{t}} = \frac{1}{L}\sum_{i=1}^{L}\hat{t}(Y_i)$$

is the estimator sample mean. We estimated the estimator bias gradient length (BGL) (the norm $\|\cdot\|_C$ with $C = I$) via the methods described in Section III. We traced out the estimator performance curves in the delta-sigma plane by varying the regularization parameter $\beta$.

Fig. 9 illustrates the benefits of using the bootstrap estimate of BGL as compared with the ordinary method-of-moments BGL estimator for the identity penalized likelihood estimator. Included are standard error bars (twice the length gives 95% confidence intervals) for bias (horizontal lines) and variance (vertical lines smaller than plotting symbol) of the bootstrap BGL estimator for $L = 500$ and $L = 1000$ realizations. The
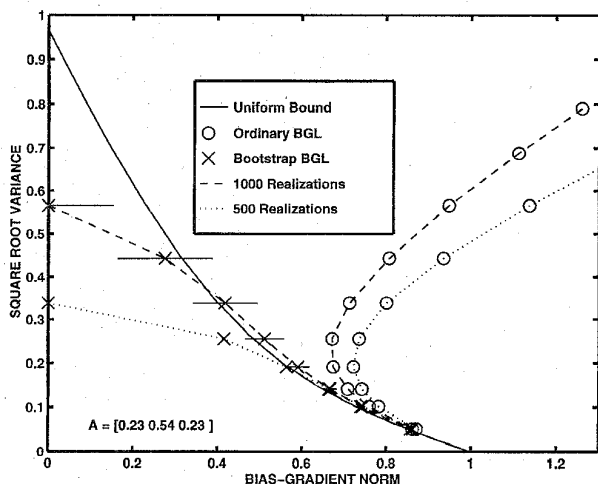
Fig. 9. Performance of penalized likelihood estimator compared to uniform CR bound. Bias gradient length (BGL) estimates were computed using both the standard method-of-moments estimate and the bootstrap estimate described in Section III. Data points to left fall below bound due to insufficient number of realizations for reliable BGL estimation.



Fig. 10. Performance of three different estimators compared with the uniform CR bound.

BGL error bars were computed under a large $L$ Gaussian approximation to the bias gradient estimates and a square root transformation. In general, as the smoothing parameter $\beta$ is decreased, QPML estimator bias decreases while QPML estimator variance increases. This increase in variance produces an increasingly large positive bias in the ordinary BGL estimator, causing the curve to abruptly diverge to the right. The bootstrap BGL estimator diverges for a much smaller value of $\beta$ and therefore extends the range of $\delta$, which can be reliably studied.

In Fig. 10, we compare the three different estimators to the uniform CR bound. As predicted by the asymptotic analysis the uniform bound is virtually achieved by the identity penalized likelihood estimator in the high bias and low variance region (large $\beta$). The identity penalized maximum likelihood estimator visibly outperforms the other two estimators. Unfortunately, for fixed $L = 1000$, as the estimator performance curves approach the left side of the delta-sigma plane, the bootstrap BGL estimates become increasingly variable (recall error bars in Fig. 9); therefore, an increasingly large number of realizations is required to make reliable comparisons between the estimator performance and the bound. On the other hand, ECT images corresponding to such highly variable estimates of $\underline{\theta}$ are unlikely to be of much practical interest.

## V. CONCLUSIONS

We have presented a method for specifying a lower bound in the delta-sigma plane defined as the set of pairs $(\delta_{\underline{\theta}}, \sigma_{\underline{\theta}})$, where $\delta_{\underline{\theta}}$ is the estimator bias gradient norm, and $\sigma^2_{\underline{\theta}}$ is the estimator variance. For two inverse problems, one linear and one nonlinear, we have established that the bound is achievable under certain circumstances.

There remain several open problems. In ill-posed problems, the Fisher matrix is singular, and an eigendecomposition appears to be required to compute the bound. For small ill-posed problems, this is not a major impediment. However,
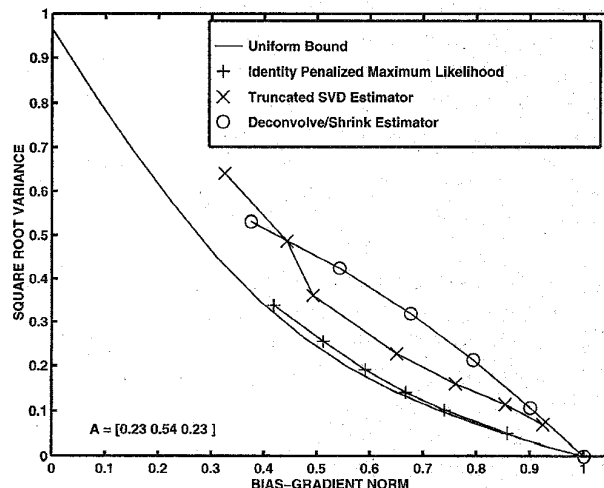
for large problems with many parameters, which includes many image reconstruction and image restoration problems, the eigendecomposition is not practical, and faster numerical methods are needed. Another problem is that the variance of the bootstrap estimator for bias gradient norm increases rapidly with the number of unknown parameters. Since the bootstrap estimator is not guaranteed to be nonnegative, this high variance can make the estimator useless for estimating small-valued bias gradient norms. In such cases, asymptotic bias and variance formulas may be useful and can be derived along similar lines as described in Appendix D. Finally, we established a general relation between bias gradient norm and maximal bias variation. Although for general estimation problems the interpretation of the bias gradient norm may be difficult, for the two applications considered in this paper, the bias gradient norm has a natural interpretation as a measure of spatial resolution of the estimator.

## APPENDIX A
## PROOF OF THEOREM 1

For a fixed $\delta > 0$, we perform constrained minimization of the biased form of the CR bound (7) over the feasible set $\nabla b_{\underline{\theta}} \colon \|\nabla b_{\underline{\theta}}\|_C \leq \delta$ of bias gradient vectors

$$\text{var}_{\underline{\theta}}(\hat{t}) \geq [\nabla t_{\underline{\theta}} + \nabla b_{\underline{\theta}}]^T F_Y^+ [\nabla t_{\underline{\theta}} + \nabla b_{\underline{\theta}}]$$

$$\geq \min_{\underline{d} \colon \|\underline{d}\|_C \leq \delta} Q(\underline{d})$$

where

$$Q(\underline{d}) = [\nabla t_{\underline{\theta}} + \underline{d}]^T F_Y^+ [\nabla t_{\underline{\theta}} + \underline{d}]$$

and $\underline{d}$ is a vector in $\mathbb{R}^n$. Defining $\tilde{\underline{d}} = C^{1/2}\underline{d}$, where $C^{1/2}$ is a square root factor of $C$, the minimization of $Q(d)$ is equivalent to

$$\min_{\tilde{\underline{d}} \colon \|\tilde{\underline{d}}\|_2 \leq \delta} [C^{1/2}\nabla t_{\underline{\theta}} + \tilde{\underline{d}}]^T B [C^{1/2}\nabla t_{\underline{\theta}} + \tilde{\underline{d}}] \qquad (45)$$

where $B = C^{-(1/2)} F_Y^+ C^{-(1/2)}$.

First, we consider the case where the unconstrained minimum $Q(\underline{d}) = 0$ occurs in the interior of the constraint set

$||\underline{d}||_C \leq \delta$. From (45), it is clear that $Q(\underline{d})$ can be zero if and only if $C^{1/2}\nabla t_\theta + \underline{d}$ lies in the null space of $B$. Such a solution $\underline{d}^o$ must have the form

$$\underline{d}^o = -\mathcal{P}_B C^{1/2}\nabla t_\theta + \underline{\phi}$$

where $\underline{\phi}$ is an arbitrary vector in the null space of $B$. However, for $\underline{d}^o$ to be a feasible solution, it must satisfy $||\underline{d}^o||_2 \leq \delta$ so that, by orthogonality of $\mathcal{P}_B C^{1/2}\nabla t_\theta$ and $\underline{\phi}$

$$\delta^2 \geq ||\underline{d}^o||_2^2 = ||\mathcal{P}_B C^{1/2}\nabla t_\theta||_2^2 + ||\underline{\phi}||_2^2 \geq ||\mathcal{P}_B C^{1/2}\nabla t_\theta||_2^2.$$

We conclude that $\min_{\underline{d}:||\underline{d}||_C \leq \delta} Q(\underline{d}) = 0$ iff $||\mathcal{P}_B C^{1/2}\nabla t_\theta||_2^2 = \nabla^T t_\theta C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_\theta \leq \delta^2$.

If $\nabla^T t_\theta C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_\theta = 0$, then we have nothing left to prove. Otherwise, assume $\delta$ lies in the range $0 \leq \delta^2 < \nabla^T t_\theta C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_\theta$. In this case, the minimizing $\underline{\tilde{d}}$ lies on the boundary and satisfies the equality constraint $||\underline{\tilde{d}}||_2 = \delta$. We thus need to solve the unconstrained minimization of the Lagrangian:

$$\min_{\underline{\tilde{d}}} [[C^{1/2}\nabla t_\theta + \underline{\tilde{d}}]^T B[C^{1/2}\nabla t_\theta + \underline{\tilde{d}}] + \lambda(\underline{\tilde{d}}^T \underline{\tilde{d}} - \delta^2)] \quad (46)$$

where we have introduced the undetermined multipler $\lambda \geq 0$. Assuming for the moment that $\lambda$ is strictly positive, the matrix $\lambda I + B$ is positive definite, and the completion of the square in the Lagrangian in (46) gives

$$[\underline{\tilde{d}} + (\lambda I + B)^{-1}BC^{1/2}\nabla t_\theta]^T (\lambda I + B)$$
$$\cdot [\underline{\tilde{d}} + (\lambda I + B)^{-1}BC^{1/2}\nabla t_\theta] + \nabla^T t_\theta F_Y^+ \nabla t_\theta$$
$$- \nabla^T t_\theta [C^{1/2}B(\lambda I + B)^{-1}BC^{1/2}]\nabla t_\theta$$
$$- \lambda\delta^2.$$

It follows immediately from the above that

$$\underline{\tilde{d}} = \underline{\tilde{d}}_{\min} = -(\lambda I + B)^{-1}BC^{1/2}\nabla t_\theta$$

achieves the minimum. Noting that $\underline{d}_{\min} = C^{-(1/2)}\underline{\tilde{d}}_{\min}$, expressing $B$ in terms of $C$ and $F_Y$, and performing simple matrix algebra, we obtain (10). Substituting the expression for $\underline{\tilde{d}}_{\min}$ into (45)

$$\min_{\underline{d}:||\underline{d}||_C \leq \delta} Q(\underline{d})$$
$$= [\nabla t_\theta + \underline{d}_{\min}]^T F_Y^+ [\nabla t_\theta + \underline{d}_{\min}]$$
$$= \nabla^T t_\theta C^{1/2}[I - [\lambda I + B]^{-1}B]^T B[I - [\lambda I + B]^{-1}B]$$
$$\cdot C^{1/2}\nabla t_\theta$$
$$= \lambda^2 \nabla^T t_\theta C^{1/2}[\lambda I + B]^{-1}B[\lambda I + B]^{-1}C^{1/2}\nabla t_\theta \quad (47)$$

which, after simple matrix manipulations, gives (9).

The Lagrange multiplier $\lambda$ is determined by the equality constraint

$$\delta^2 = \underline{\tilde{d}}_{\min}^T \underline{\tilde{d}}_{\min} = \nabla^T t_\theta [C^{1/2}B[\lambda I + B]^{-2}BC^{1/2}]\nabla t_\theta$$
$$= g(\lambda). \quad (48)$$

Substitution of $B = C^{-(1/2)}F_Y^+ C^{-(1/2)}$ yields (11) after simple matrix algebra.

Let the nonnegative definite symmetric matrix $B$ have eigendecomposition $B = \Sigma_{i=1}^r \beta_i \underline{\xi}_i \underline{\xi}_i^T$, where

$\beta_i$    positive eigenvalues;
$\underline{\xi}_i$    eigenvectors;
$r > 0$ rank of $B$.

With these definitions, the function $g(\lambda)$ (48) has the equivalent form

$$g(\lambda) = \sum_{i=1}^r \frac{\beta_i^2}{(\lambda + \beta_i)^2}|\nabla^T t_\theta C^{1/2}\underline{\xi}_i|^2. \quad (49)$$

Since by assumption $\nabla^T t_\theta C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_\theta > 0, C^{1/2}\nabla t_\theta$ does not lie in the nullspace of $B$, and thus, $|\nabla^T t_\theta C^{1/2}\underline{\xi}_i|^2 > 0$ for at least one $i, i = 1, \cdots, r$. Therefore, from (49), it is obvious that the function $g(\lambda)$ is continuous monotone decreasing over $\lambda \geq 0$ with $\lim_{\lambda\to\infty} g(\lambda) = 0$, and $\lim_{\lambda\to 0} g(\lambda) = \nabla t_\theta C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_\theta$. Hence, there exists a unique strictly positive $\lambda$ such that $g(\lambda) = \delta^2$ for any value $\delta^2 \in [0, \nabla t_\theta C^{1/2}\mathcal{P}_B C^{1/2}\nabla t_\theta)$. $\square$

## APPENDIX B
### BOOTSTRAP DERIVATION

We start with the following simple estimator $\widehat{\nabla b_\theta} = \widehat{\nabla b_\theta}(Y_1, \cdots, Y_L)$ of the bias gradient $\nabla b_\theta$

$$\widehat{\nabla b_\theta} = \frac{1}{L}\sum_{i=1}^L \underline{z}(Y_i)$$

where $\underline{z}$ is the column vector

$$\underline{z}(Y_i) = \hat{t}(Y_i)\nabla_\theta^T \ln f_Y(Y_i; \underline{\theta}) - \nabla t_\theta.$$

The (biased) sample magnitude estimator of the norm squared $\delta^2 = ||\nabla b_\theta||_C^2$ is

$$\hat{\delta}^2(Y_1, \cdots, Y_L) = ||\widehat{\nabla b_\theta}||_C^2 = \frac{1}{L^2}\left\|\sum_{i=1}^L \underline{z}(Y_i)\right\|_C^2.$$

Now, given the random sample $Y_1 = y_1, \cdots, Y_L = y_L$, the resampled estimate $\hat{\delta}_*^2 = \hat{\delta}_*^2(Y_1^*, \cdots, Y_L^*)$ is [37]

$$\hat{\delta}_*^2 = \frac{1}{L^2}\left\|\sum_{i=1}^L c_i\underline{z}(y_i)\right\|_C^2$$
$$= \frac{1}{L^2}\sum_{i=1}^L\sum_{j=1}^L c_ic_j\langle\underline{z}(y_i), \underline{z}(y_j)\rangle_C \quad (50)$$

where $\langle\underline{u}, \underline{v}\rangle_C = \underline{u}^T C\underline{v}$ is defined as the (weighted) inner product of column vectors $\underline{u}$ and $\underline{v}$. Define $\underline{c} = [c_1, \cdots, c_L]^T$, and let $H = [[\langle\underline{z}(Y_i), \underline{z}(Y_j)\rangle_C]]$ denote a $L \times L$ matrix of inner products. Then, the resampled estimate (50) has the equivalent form

$$\hat{\delta}_*^2 = \frac{1}{L^2}\underline{c}^T H\underline{c} \quad (51)$$
$$= \frac{1}{L^2}\text{trace}\{H\underline{cc}^T\}. \quad (52)$$

The resampling outcomes $c_1, \cdots, c_L$ are multinomial distributed with equal cell probabilities $p_1 = \cdots = p_L = 1/L$ and $\Sigma_{i=1}^L c_i = L$. Averaging (52) over $\underline{c}$, we obtain the bootstrap estimate of the mean

$$E_*[\hat{\delta}_*^2] = \sum_* \hat{\delta}_*^2 \binom{L}{c_1 \cdots c_L} L^{-L}$$
$$= \frac{1}{L^2}\text{trace}\{HE_*[\underline{cc}^T]\}. \quad (53)$$

From the mean and covariance of the multinomial distribution [49, Sec. 3.2]:

$$E_*[\underline{c}\,\underline{c}^T] = I + \frac{L-1}{L}\underline{1}\underline{1}^T$$

where $\underline{1}$ is a vector of ones. Applying this identity to (53) and noting that trace$\{H\} = \Sigma_{i=1}^L \|\underline{z}(Y_i)\|_C^2$ and $\underline{1}^T H \underline{1} = \|\Sigma_{i=1}^L \underline{z}(Y_i)\|_C^2$, we obtain

$$E_*[\hat{\delta}^2*] = \frac{1}{L^2}\text{trace}\,\{HE_*[\underline{c}\,\underline{c}^T]\}$$

$$= \frac{1}{L^2}\left[\text{trace}\,\{H\} + \frac{L-1}{L}\underline{1}^T H\underline{1}\right]$$

$$= \frac{1}{L^2}\left[\sum_{i=1}^L \|\underline{z}(Y_i)\|_C^2 - L\|\bar{z}\|_C^2 + L^2\|\bar{z}\|_C^2\right]$$

$$= \frac{1}{L^2}\sum_{i=1}^L \|\underline{z}(Y_i) - \bar{z}\|_C^2 + \hat{\delta}^2$$

where

$$\bar{z} = \overline{\nabla b_\theta} = \frac{1}{L}\sum_{i=1}^L \underline{z}(Y_i)$$

is the sample mean, and we have identified $\hat{\delta}^2 = \|\bar{z}\|_C^2$. Plugging this last expression into (26), we obtain

$$\hat{\delta}_c^2 = 2\hat{\delta}^2 - E_*[\hat{\delta}_*^2]$$

$$= \hat{\delta}^2 - \frac{1}{L^2}\sum_{i=1}^L \|\underline{z}(Y_i) - \bar{z}\|_C^2$$

which is identical to (27).

## VIII. APPENDIX C
### LOWER BOUND ON BIAS GRADIENT

Here, we derive a simple lower bound on the maximal bias variation over the region $\mathcal{C} = \{\underline{u}: (\underline{u} - \underline{\theta})^T C^{-1}(\underline{u} - \underline{\theta}) \leq 1\}$ under the assumptions i) $F_Y(\underline{u})$ is constant over $\underline{u} \in \mathcal{C}$, and ii) the functional $t_{\underline{u}}$ to be estimated is linear over $\underline{u} \in \mathcal{C}$. Define

1) $\mathcal{P}_F = F_Y F_Y^+ = F_Y^+ F_Y$ as the symmetric matrix that orthogonally projects onto the range space of $F_Y$,
2) $\mathcal{N}_F = \{\underline{u}: \mathcal{P}_F \underline{u} = \underline{0}\}$ as the nullspace of $F_Y$
3) $m_{\underline{u}} = E_{\underline{u}}[t]$.

Under assumption i), the parameter $\underline{u}$ is not identifiable for $\underline{u} \in \mathcal{N}_F \cap \mathcal{C}$, and it follows that $m_{\underline{u}} - m_{\underline{\theta}} = 0$. Therefore, we obtain the lower bound

$$\max_{\underline{u}\in\mathcal{C}} |b_{\underline{u}} - b_{\underline{\theta}}|^2 = \max_{\underline{u}\in\mathcal{C}} |m_{\underline{u}} - t_{\underline{\theta}} - m_{\underline{\theta}} + t_{\underline{\theta}}|^2$$

$$\geq \max_{\underline{u}\in\mathcal{C}\cap\mathcal{N}_F} |m_{\underline{u}} - m_{\underline{\theta}} + t_{\underline{\theta}} - t_{\underline{u}}|^2$$

$$= \max_{\underline{u}\in\mathcal{C}\cap\mathcal{N}_F} |t_{\underline{u}} - t_{\underline{\theta}}|^2$$

$$= \max_{\underline{u}\in\mathcal{C}\cap\mathcal{N}_F} |\underline{h}^T \Delta\underline{u}|^2$$

$$= \max_{\underline{u}\in\mathcal{C}} |\underline{h}^T[I - \mathcal{P}_F]\Delta\underline{u}|^2 \qquad (54)$$

where $\Delta\underline{u} = \underline{u} - \underline{\theta}$. Now, using an extremal property of the Rayleigh quotient, the right-hand side of (54) is

$$\max_{\underline{u}\in\mathcal{C}} |\underline{h}^T[I - \mathcal{P}_F]\Delta\underline{u}|^2$$

$$= \max_{\Delta\underline{u}} \frac{\Delta\underline{u}^T[I - \mathcal{P}_F]\underline{h}\underline{h}^T[I - \mathcal{P}_F]\Delta\underline{u}}{\Delta\underline{u}^T C^{-1}\Delta\underline{u}}$$

$$= \max_{\tilde{u}} \frac{\tilde{u}^T C^{1/2}[I - \mathcal{P}_F]\underline{h}\underline{h}^T[I - \mathcal{P}_F]C^{1/2}\tilde{u}}{\tilde{u}^T \tilde{u}}$$

$$= \underline{h}^T[I - \mathcal{P}_F]C[I - \mathcal{P}_F]\underline{h} \qquad (55)$$

where, in the third line, $\tilde{u} = C^{-(1/2)}\Delta\underline{u}$.

In view of (6), the combination of (54) and (55) yields the following lower bound on the norm:$\|\nabla b_{\underline{u}}\|_C$:

$$\|\nabla b_{\underline{u}}\|_C \geq \sqrt{\underline{h}^T[I - \mathcal{P}_F]C[I - \mathcal{P}_F]\underline{h}} + \epsilon. \qquad (56)$$

For the case that the bias $b_{\underline{u}}$ is linear over $\underline{u} \in \mathcal{C}$, $\epsilon$ in (56) is equal to zero, and we have an exact bound $\|\nabla b_{\underline{u}}\|_C \geq \sqrt{\underline{h}^T[I - \mathcal{P}_F]C[I - \mathcal{P}_F]\underline{h}}$. In view of the relation (54), replacing $\underline{h}$ with $\nabla t_\theta$ in this bound will give an approximate bound when $t_{\underline{u}}$ is nonlinear but smooth over $\underline{u} \in \mathcal{C}$. Expression (56) with $\epsilon = 0$ will probably be a fairly good approximation to the bound when the range space components of $\underline{u}$ can be estimated without bias. This is true for linear models $Y = A\underline{u} + \underline{w}$. However, the reader is cautioned that for nonlinear models, (56) may not be very tight since unbiased estimators may not exist even for components lying in the range space of $F_Y$ [42].

### APPENDIX D
### ASYMPTOTIC APPROXIMATION OF BIAS, BIAS GRADIENT, AND VARIANCE FOR POISSON QPML

Define the vector

$$\underline{z} = [F_Y(\underline{\theta}) + \beta P]^{-1} F_Y(\underline{\theta})\underline{\theta},$$

$$= \left[I - \left[P + \frac{1}{\beta}F_Y(\underline{\theta})\right]^{-1}P\right]\underline{\theta}. \qquad (57)$$

Here, we derive the following asymptotic formulas for variance, bias, and bias gradient of the Poisson QPML estimator of a general differentiable function $t_{\underline{\theta}}$.

**Asymptotic Variance:**

$$\sigma_{\underline{\theta}}^2 = \nabla^T t_{\underline{z}}[F_Y(\underline{\theta}) + \beta P]^{-1}F_Y(\underline{\theta})[F_Y(\underline{\theta}) + \beta P]^{-1}\nabla t_{\underline{z}}. \qquad (58)$$

**Asymptotic Bias:**

$$b_{\underline{\theta}} = t_{\underline{z}} - t_{\underline{\theta}} \qquad (59)$$

**Asymptotic Bias Gradient:**

$$\nabla b_{\underline{\theta}} = F_Y(\underline{\theta})[\beta P + F_Y(\underline{\theta})]^{-1}\nabla t_{\underline{z}} - \nabla t_{\underline{\theta}} - O\left(\frac{1}{\beta}\right) \qquad (60)$$

where

$$O\left(\frac{1}{\beta}\right) = \frac{1}{\beta}\sum_{j=1}^m \gamma_j(\underline{\theta})A_{j*}^T B^{-1}\nabla t_{\underline{z}} \qquad (61)$$

$$\gamma_j(\underline{\theta}) = \frac{A_{j*}\underline{\theta}^T PB^{-1}A_{j*}}{\mu_j^2(\underline{\theta})} \qquad (62)$$

and $\nabla t_{\underline{z}}$ denotes the evaluation of the gradient of $t_{\underline{\theta}}$ at the point $\underline{\theta} = \underline{z}$.

Define the ambiguity function $a(\underline{u}, \underline{\theta}) = E_{\underline{\theta}}[J(\underline{u})]$, and let $\underline{u} = \underline{z} = \underline{z}(\underline{\theta})$ be the root of the equation $\underline{\lambda}(\underline{u}) = 0$, where $\underline{\lambda}(\underline{u}) = \nabla^{10} a(\underline{u}, \underline{\theta})$. Assuming the technical conditions underlying [50, Corollary 3.2, Sec. 6.3] are satisfied,[1] we have the following approximation: In the limit of large observation time, the estimator $\hat{\underline{\theta}}$ is asymptotically normal with mean $\underline{z}$ and covariance matrix $\Sigma = [\nabla^{20} a(\underline{z}, \underline{\theta})]^{-1} G_{\underline{z}} [\nabla^{20} a(\underline{z}, \underline{\theta})]^{-T}$, where $G_{\underline{z}} = \text{cov}_{\underline{\theta}}(\nabla J(\underline{z}))$. Furthermore, assuming that the function $t_{\underline{\theta}}$ has nonzero derivative at $\underline{\theta} = \underline{z}$ the estimator $\hat{t} = t_{\hat{\underline{\theta}}}$ is asymptotically normal with mean $t_{\underline{z}}$ and variance $\nabla^T t_{\underline{z}} \Sigma \nabla t_{\underline{z}}$ [51, p. 122, Theorem A]. This gives the asymptotic expression for bias $b_{\underline{\theta}} = E_{\underline{\theta}}[\hat{t} - t_{\underline{\theta}}] = t_{\underline{z}} - t_{\underline{\theta}}$ and an asymptotic expression for variance

$$\text{var}_{\underline{\theta}}(\hat{t}) = \nabla^T t_{\underline{z}} [\nabla^{20} a(\underline{z}, \underline{\theta})]^{-1} G_{\underline{z}} [\nabla^{20} a(\underline{z}, \underline{\theta})]^{-1} \nabla t_{\underline{z}}. \quad (63)$$

Since the penalized Poisson likelihood function $J(\underline{u})$ in (42) is linear in the observations $\underline{Y}$, and $Y_i$ are independent Poisson random variables, it is simple to derive the following expression for the covariance matrix of $\nabla J(\underline{z}) = \sum_{j=1}^T A_{j*} (Y_j / \mu_j(\underline{z}) - 1) - \beta P \underline{z}$:

$$\begin{aligned} \text{cov}_{\underline{\theta}}[\nabla J(\underline{z})] &= F(\underline{z}, \underline{\theta}) \\ &= \sum_{j=1} A_{j*} A_{j*}^T \frac{1}{\mu_j(\underline{\theta})} \left( \frac{\mu_j(\underline{\theta})}{\mu_j(\underline{z})} \right)^2 \\ &= F_Y(\underline{\theta}) + o(\mu(\underline{z} - \underline{\theta})) \end{aligned} \quad (64)$$

where $A_{j*}^T$ is the $j$th row of $A, F_Y(\underline{\theta})$ is the Fisher matrix (29), and $o(\mu(\underline{z} - \underline{\theta}))$ is a remainder term that is close to zero when the difference between the projections $\underline{\mu}(\underline{z})$ and $\underline{\mu}(\underline{\theta})$ is small. To obtain (64) with remainder term, we used the series development

$$\frac{\mu_j(\underline{\theta})}{\mu_j(\underline{z})} = 1 - \frac{1}{\mu_j(\underline{\theta})} (\underline{z} - \underline{\theta})^T A_{j*} + o((\underline{z} - \underline{\theta})^T A_{j*}).$$

The ambiguity function is

$$a(\underline{u}, \underline{\theta}) = \sum_{j=1}^m (\mu_j(\underline{\theta}) \ln \mu_j(\underline{u}) - \mu_j(\underline{u})) - \frac{\beta}{2} \underline{u}^T P \underline{u}.$$

Differentiation of the ambiguity function with respect to $\underline{u}$ yields

$$\begin{aligned} \nabla^{10} a(\underline{u}, \underline{\theta}) &= \sum_{j=1}^m A_{j*} \left( \frac{\mu_j(\underline{\theta})}{\mu_j(\underline{u})} - 1 \right) - \beta P \underline{u} \\ &= - \left( \sum_{j=1}^m \frac{1}{\mu_j(\underline{\theta})} A_{j*}^T A_{j*} \right) (\underline{u} - \underline{\theta}) - \beta P \underline{u} \\ &\quad + o(\mu(\underline{u} - \underline{\theta})) \\ &= -F_Y(\underline{\theta})(\underline{u} - \underline{\theta}) - \beta P \underline{u} + o(\mu(\underline{u} - \underline{\theta})), \quad (65) \end{aligned}$$

and similarly,

$$\nabla^{20} a(\underline{u}, \underline{\theta}) = -F_Y(\underline{\theta}) - \beta P + O(\mu(\underline{u} - \underline{\theta})). \quad (66)$$

[1] Among other things, these conditions involve showing that the gradient function $\nabla J(\underline{\theta})$ converges a.s. to zero as the observation time increases.

Neglecting the $O(\mu(\underline{u} - \underline{\theta}))$ remainder terms, multiplication of the inverse of (66) the covariance (64) and the inverse transpose of (66) yields the asymptotic variance expression (63). Likewise, neglecting the remainder term in (65) and solving for the root $\underline{u} = \underline{z}$ of the equation $\nabla^{10} a(\underline{u}, \underline{\theta}) = 0$ yields the asymptotic expression for the root (57).

We next derive (60) for the bias gradient. Applying the chain rule to differentiate (59), we obtain

$$\nabla b_{\underline{\theta}} = \nabla \underline{z}_{\underline{\theta}}^T \nabla t_{\underline{z}} - \nabla t_{\underline{\theta}} \quad (67)$$

where $\nabla \underline{z}_{\underline{\theta}}^T$ is an $n \times n$ matrix of derivatives of $\underline{z} = \underline{z}_{\underline{\theta}}$. From (57), the $k$th row of this matrix is

$$\begin{aligned} \frac{d}{d\theta_k} \underline{z}_{\underline{\theta}}^T &= \underline{e}_k^T \left[ I - P \left[ P + \frac{1}{\beta} F_Y(\underline{\theta}) \right]^{-1} \right] \\ &\quad - \underline{\theta}^T P \frac{d}{d\theta_k} \left[ P + \frac{1}{\beta} F_Y(\underline{\theta}) \right]^{-1} + O(\mu(\underline{u} - \underline{\theta})). \end{aligned} \quad (68)$$

Define the matrix $B = P + (1/\beta) F_Y(\underline{\theta})$. From the differentiation formula $(d/d\theta_k) B(\theta)^{-1} = -B^{-1}(d/d\theta_k)B(\theta)B^{-1}$ and from (41) for the Fisher information matrix $F_Y(\underline{\theta})$, we have

$$\begin{aligned} -\underline{\theta}^T P \frac{d}{d\theta_k} \left[ P + \frac{1}{\beta} F_Y(\underline{\theta}) \right]^{-1} &= \frac{1}{\beta} \underline{\theta}^T P B^{-1} \frac{d}{d\theta_k} F_Y(\underline{\theta}) B^{-1} \\ &= -\frac{1}{\beta} \sum_{j=1}^m \frac{A_{jk} \underline{\theta}^T P B^{-1} A_{j*}}{\mu_j^2(\underline{\theta})} A_{j*}^T B^{-1} \\ &= -\frac{1}{\beta} \sum_{j=1}^m \gamma_{jk}(\underline{\theta}) A_{j*}^T B^{-1} \end{aligned} \quad (69)$$

where $\gamma_{jk}(\underline{\theta})$ is the $k$th element of the vector $\underline{\gamma}_j$ defined in (62). Combining (68) and (69), we obtain

$$\nabla \underline{z}_{\underline{\theta}} = \left[ I - P \left[ P + \frac{1}{\beta} F_Y(\underline{\theta}) \right]^{-1} \right] + O\left(\frac{1}{\beta}\right)$$

which, when substituted into (67) and neglecting the remainder term $O(\mu(\underline{u} - \underline{\theta}))$, yields the bias gradient expression (60).
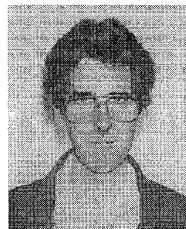
## REFERENCES

[1] S. M. Kay, Modern Spectral Estimation: Theory and Application. Englewood Cliffs, NJ: Prentice Hall, 1988.
[2] P. J. Daniell, "Discussion of paper by M. S. Bartlett," J. Royal Stat. Soc., Ser. B, vol. 8, p. 27, 1946.
[3] R. B. Blackman and J. W. Tukey, "The measurement of power spectra from the point of view of communication engineering," Bell Syst. Tech. J., vol. 37, pp. 183–282, 1958.
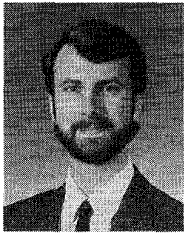[4] R. Lagendijk and J. Biemond, Iterative Identification and Restoration of Images. Boston: Kluwer, 1991.

[5] C. M. Stein, "Multiple regression," in *Contributions to Probability and Statistics: Essays in Honor of Harold Hotelling*, I. Olkin, S. Ghurge, W. Hoeffding, W. Nadow, and H. Mann, Eds. Stanford, CA: Stanford Univ. Press, 1960, pp. 424–443.

[6] A. E. Hoerl and R. W. Kennard, "Ridge regression: Biased estimation of nonorthogonal components," *Technometrics*, vol. 12, pp. 55–67, 1970.

[7] J.-S. Liow and S. C. Strother, "Practical tradeoffs between noise, quantitation, and number of iterations for maximum likelihood-based reconstructions," *IEEE Trans. Med. Imaging*, vol. 10, no. 4, pp. 563–571, 1991.

[8] S. J. Lee, G. R. Gindi, I. G. Zubal, and A. Rangarajan, "Using ground-truth data to design priors in Bayesian SPECT reconstruction," in *Information Processing in Medical Imaging*, Y. Bizais, C. Barillot, and R. D. Paola, Eds. Boston, MA: Kluwer, 1995.

[9] J. A. Fessler, "Penalized weighted least-squares image reconstruction for positron emission tomography," *IEEE Trans. Med. Imaging*, vol. 13, pp. 290–300, June 1994.

[10] R. E. Carson, Y. Yan, B. Chodkowski, T. K. Yap, and M. E. Daube-Witherspoon, "Precision and accuracy of regional radioactivity quantitation using the maximum likelihood EM reconstruction algorithm," *IEEE Trans. Med. Imaging*, vol. 13, pp. 526–537, Sept. 1994.

[11] D. Wang, V. Haese-Coat, A. Bruno, and J. Ronsin, "Some statistical properties of mathematical morphology," *IEEE Trans. Signal Processing*, vol. 43, pp. 1955–1965, Aug. 1995.

[12] N. Himayat and S. A. Kassam, "Approximate performance analysis of edge preserving filters," *IEEE Trans. Signal Processing*, vol. 41, pp. 2764–2777, Sept. 1993.

[13] M. B. Woodroofe and J. W. Van Ness, "The maximum deviation of sample spectral densities," *Ann. Math. Stat.*, vol. 38, pp. 1558–1569, 1967.

[14] P. Bloomfield, *Fourier Analysis of Time Series*. New York: Wiley, 1976.

[15] A. O. Hero, "A Cramer-Rao type lower bound for essentially unbiased parameter estimation," MIT Lincoln Lab., Lexington, MA, Jan. 1992, Tech. Rep. 890, DTIC AD-A246666.

[16] J. Fessler and A. Hero, "Cramer–Rao bounds for biased estimators in image restoration," in *Proc. 36th IEEE Midwest Symp. Circuits Syst.*, Detroit, MI, Aug. 1993.

[17] M. Usman, A. Hero, J. A. Fessler, and W. Rogers, "Bias-variance tradeoffs analysis using uniform CR bound for a SPECT system," in *Proc. IEEE Nuclear Sci. Symp. Med. Imaging Conf.*, San Francisco, CA, Nov. 1993, pp. 1463–1467.

[18] M. Usman, A. Hero, and J. A. Fessler, "Bias-variance tradeoffs analysis using uniform CR bound for images," in *Proc. IEEE Image Processing Conf.*, Austin, TX, Nov. 1994, pp. 835–839.

[19] C. R. Rao, *Linear Statistical Inference and Its Applications*. New York: Wiley, 1973.

[20] J. A. Fessler and A. O. Hero, "Space-alternating generalized EM algorithm," *IEEE Trans. Signal Processing*, vol. 42, no. 10, pp. 2664–2677, Oct. 1994.

[21] G. E. P. Box and N. R. Draper, "A basis for the selection of a response surface design," *J. Am. Statist. Assoc.*, vol. 54, pp. 622–654, 1959.

[22] M. A. Ali, "On a class of shrinkage estimators of the vector of shrinkage coefficients," *Commun. Statist.-Theory Meth.*, vol. 18, no. 12, pp. 4491–4500, 1981.

[23] T. Bednarski, "On minimum bias and variance estimation for parametric models with shrinking contamination," *Prob. Math. Stat.*, vol. 6, no. 2, pp. 121–129, 1985.

[24] M. Quenouille, "Approximate tests of correlation in time series," *J. Royal Stat. Soc., Ser. B*, vol. 11, pp. 18–84, 1949.

[25] J. Tukey, "Bias and confidence in not quite large samples," *Ann. Math. Stat.*, vol. 29, pp. 614, 1956.

[26] B. Efron, "Bootstrap methods: Another look at the jackknife," *Ann. Stat.*, vol. 7, pp. 1–26, 1979.

[27] I. A. Ibragimov and R. Z. Has'minskii, *Statistical Estimation: Asymptotic Theory*. New York: Springer-Verlag, 1981.

[28] M. Fréchet, "Sur l'extension de certaines évaluations statistiques de petits echantillons," *Rev. Inst. Int. Stat.*, vol. 11, pp. 128–205, 1943.

[29] G. Darmois, "Sur les lois limites de la dispersion de certaines estimations," *Rev. Inst. Int. Stat.*, vol. 13, pp. 9–15, 1945.

[30] H. Cramér, "A contribution to the theory of statistical estimation," *Skand. Aktuaries Tidskrift*, vol. 29, pp. 458–463, 1946.

[31] C. R. Rao, "Minimum variance and the estimation of several parameters," *Proc. Cambridge Phil. Soc.*, vol. 43, pp. 280–283, 1946.

[32] H. L. Van-Trees, *Detection, Estimation, and Modulation Theory: Part I*. New York: Wiley, 1968.

[33] J. D. Gorman and A. O. Hero, "Lower bounds for parametric estimation with constraints," *IEEE Trans. Inform. Theory*, vol. 36, pp. 1285–1301, Nov. 1990.

[34] M. Usman, "Biased and unbiased Cramer-Rao bounds: Computational issues and applications," Ph.D. Thesis, Dept. Elec. Eng. Comput. Sci., Univ. Michigan, Aug. 1994.

[35] W. Press, S. Teukolsky, W. Vetterling, and B. Flannery, *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge, UK: Cambridge Univ. Press, 1992, 2nd ed.

[36] A. O. Hero, M. Usman, A. C. Sauve, and J. A. Fessler, "Recursive algorithms for computing the Cramer-Rao bound," *IEEE Trans. Signal Processing*, submitted for publication.

[37] B. E. Efron, *The Jackknife, the Bootstrap, and Other Resampling Plans*. Philadelphia, PA: SIAM, 1982.

[38] V. A. Morozov, *Methods for Solving Incorrectly Posed Problems*. New York: Springer-Verlag, 1984.

[39] K. Lange, "Convergence of EM image reconstruction algorithms with Gibbs smoothing," *IEEE Trans. Med. Imaging*, vol. 9, pp. 439–446, June 1991.

[40] T. Hebert and R. Leahy, "A generalized EM algorithm for 3-D Bayesian reconstruction from Poisson data using Gibbs priors," *IEEE Trans. Med. Imaging*, vol. 8, pp. 194–203, 1989.

[41] V. Johnson, W. Wong, X. Hu, and C. Chen, "Image restoration using Gibbs priors: boundary modeling, treatment of blurring, and selection of hyperparameters," *IEEE Trans. Patt. Anal. Machine Intell.*, vol. PAMI-13, pp. 413–425, 1991.

[42] R. C. Liu and L. D. Brown, "Nonexistence of informative unbiased estimators in singular problems," *Ann. Stat.*, vol. 21, no. 1, pp. 1–13, 1993.

[43] G. H. Golub and C. F. Van Loan, *Matrix Computations*. Baltimore, MA: The Johns Hopkins Univ. Press, 1989, 2nd ed.

[44] D. L. Snyder and M. I. Miller, *Random Point Processes in Time and Space*. New York: Springer-Verlag, 1991.

[45] A. O. Hero and J. A. Fessler, "A recursive algorithm for computing CR-type bounds on estimator covariance," *IEEE Trans. Inform. Theory*, vol. 40, pp. 1205–1210, July 1994.

[46] J. A. Fessler and A. O. Hero, "Penalized maximum likelihood image reconstruction using space alternating generalized EM algorithms," *IEEE Trans. Image Processing*, vol. 4, pp. 1417–1429, Oct. 1995.

[47] J. A. Fessler and W. L. Rogers, "Uniform quadratic penalties cause nonuniform image resolution (and sometimes vice versa)," in *Proc. IEEE Nuclear Sci. Symp.*, vol. 4, 1994, pp. 1915–1919.

[48] J. A. Fessler, "Resolution properties of regularized image reconstruction methods," Technical Report 297, Comm. Signal Processing Lab. (CSPL), Dept. Elec. Eng. Comput. Sci., Univ. of Michigan, Ann Arbor, Aug. 1995.

[49] E. B. Manoukian, *Modern Concepts and Theorems of Mathematical Statistics*. New York: Springer-Verlag, 1986.

[50] P. J. Huber, *Robust Statistics*. New York: Wiley, 1981.

[51] R. J. Serfling, *Approximation Theorems of Mathematical Statistics*. New York: Wiley, 1980.

**Alfred O. Hero, III** (M'84) was born in Boston, MA, in 1955. He received the B.S. degree *summa cum laude* from Boston University, Boston, MA, in 1980 and the Ph.D degree from Princeton University, Princeton, NJ, in 1984, both in electrical engineering.
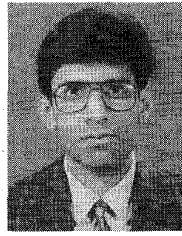
He held the G.V.N. Lothrop Fellowship in Engineering at Princeton University. He is presently Professor of Electrical Engineering and Computer Science and Area Chair for Signal Processing at the University of Michigan, Ann Arbor. He has held positions of Visiting Scientist at Lincoln Laboratory of the Massachusetts Institute of Technology, Lexington, MA, from 1987 to 1989, Visiting Professor at Ecole Nationale de Techniques Avancées (ENSTA), Paris, France, in 1991, and William Clay Ford Fellow at Ford Motor Company, Dearborn, MI, in 1993. His research interests are in the areas of detection and estimation theory applied to statistical signal and image processing.

Dr. Hero is a member of Tau Beta Pi, the American Statistical Association, and Commission C of the International Union of Radio Science (URSI). He is currently an Associate Editor for the IEEE TRANSACTIONS ON INFORMATION THEORY. He was Chairman for Publicity for the *1986 IEEE International Symposium on Information Theory*. He was General Chairman for the *1995 IEEE International Conference on Acoustics, Speech, and Signal Processing*.

**Jeffrey A. Fessler** (M'90) received the B.S.E.E. degree from Purdue University, West Lafayette, IN, in 1985, the M.S.E.E. degree from Stanford University, Stanford, CA, in 1986, and the M.S. degree in statistics from Stanford University in 1989. From 1985 to 1988, he was a National Science Foundation Graduate Fellow at Stanford, where he received the Ph.D. degree in electrical engineering in 1990.

He has worked at the University of Michigan, Ann Arbor, since then. From 1991 to 1992, he was a Department of Energy Alexander Hollaender Post-Doctoral Fellow in the Division of Nuclear Medicine. From 1993 to 1995, he was an Assistant Professor in Nuclear Medicine and the Bioengineering Program. Since 1995, he has been with the Department of Electrical Engineering and Computer Science. His research interests are in statistical aspects of imaging problems.

**Mohammad Usman** (M'95) received the B.S. degree in electrical engineering from the University of Engineering and Technology, Lahore, Pakistan, in 1987, the M.S. degree in electronics engineering from Philips International Institute, Eindhoven, the Netherlands, in 1989, and the Ph.D. degree in electrical engineering and computer science from the University of Michigan, Ann Arbor, in 1994.

From 1988 to 1989, he was at the Philips Research Laboratories, Eindhoven, the Netherlands. From 1989 to 1990, he was at the Advanced Development Center, Philips Industrial Activities, Leuven, Belgium. Since 1995, he has been on the faculty of the FAST Institute of Computer Science, Lahore, Pakistan, where he is an Assistant Professor. He is also a consultant to Cressoft, Denver, CO. His current research interests include image reconstruction, speech compression, and code optimization.